

# THE PROBABILITY OF BACKTEST OVERFITTING

David H. Bailey <sup>\*</sup>      Jonathan M. Borwein <sup>†</sup>

Marcos López de Prado <sup>‡</sup>      Qiji Jim Zhu <sup>§</sup>

February 12, 2014

First version: August 2013

---

<sup>\*</sup>Lawrence Berkeley National Laboratory (retired), 1 Cyclotron Road, Berkeley, CA 94720, USA, and Research Fellow at the University of California, Davis, Department of Computer Science. E-mail: [david@davidhbailey.com](mailto:david@davidhbailey.com); URL: <http://www.davidhbailey.com>

<sup>†</sup>Laureate Professor of Mathematics at University of Newcastle, Callaghan NSW 2308, Australia, and a Fellow of the Royal Society of Canada, the Australian Academy of Science and the AAAS. E-mail: [jonathan.borwein@newcastle.edu.au](mailto:jonathan.borwein@newcastle.edu.au); URL: <http://www.carma.newcastle.edu.au/jon>

<sup>‡</sup>Senior Managing Director at Guggenheim Partners, New York, NY 10017, and Research Affiliate at Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA. E-mail: [lopezdeprado@lbl.gov](mailto:lopezdeprado@lbl.gov); URL: <http://www.QuantResearch.info>

<sup>§</sup>Professor, Department of Mathematics, Western Michigan University, Kalamazoo, MI 49008, USA. Email: [zhu@wmich.edu](mailto:zhu@wmich.edu); URL: <http://homepages.wmich.edu/~zhu/>

# THE PROBABILITY OF BACKTEST OVERFITTING

## Abstract

Most firms and portfolio managers rely on backtests (or historical simulations of performance) to select investment strategies and allocate them capital. Standard statistical techniques designed to prevent regression overfitting, such as hold-out, tend to be unreliable and inaccurate in the context of investment backtests. We develop a framework that estimates the probability of backtest overfitting (PBO) specifically in the context of investment simulations, through a non-parametric numerical method that we call combinatorially symmetric cross-validation (CSCV). We show that CSCV produces accurate estimates of the probability that a particular backtest is overfit.

**Keywords.** Backtest, historical simulation, probability of backtest overfitting, investment strategy, optimization, Sharpe ratio, minimum backtest length, performance degradation.

**JEL Classification:** G0, G1, G2, G15, G24, E44.

**AMS Classification:** 91G10, 91G60, 91G70, 62C, 60E.

**Acknowledgements.** We are indebted to the Editor and two anonymous referees who peer-reviewed a related article for the *Notices of the American Mathematical Society* [1]. We are also grateful to Tony Anagnostakis (Moore Capital), Marco Avellaneda (Courant Institute, NYU), Peter Carr (Morgan Stanley, NYU), Paul Embrechts (ETH Zürich), Matthew D. Foreman (University of California, Irvine), Ross Garon (SAC Capital), Paul Glasserman (Columbia University), Jeffrey Lange (Guggenheim Partners), Attilio Meucci (KKR, NYU), Natalia Nolde (University of British Columbia and ETH Zürich) and Riccardo Rebonato (PIMCO, University of Oxford) for many useful and stimulating exchanges.

**Sponsorship.** Research supported in part by the Director, Office of Computational and Technology Research, Division of Mathematical, Information, and Computational Sciences of the U.S. Department of Energy, under contract number DE-AC02-05CH11231 and by various Australian Research Council grants.

“This was our paradox: No course of action could be determined by a rule, because every course of action can be made to accord with the rule.”

Ludwig Wittgenstein [26]

## 1 INTRODUCTION

Despite of its limitations, the *Sharpe ratio* (SR)<sup>1</sup> is the “gold standard” of investment performance evaluation. Bailey and López de Prado [2] developed methodologies to assess the probability that a SR is inflated (PSR), and the minimum track record length (MinTRL) required for a SR to be statistically significant. These statistics were developed to judge the reliability of SR computed on live performance (track record). We have not been able to find similar statistics or methodologies applicable to judging backtested SR. Thus began our quest for a general methodology to assess the reliability of a backtest.

**A common approach.** Perhaps the most common approach among practitioners is to require the researcher to “hold-out” a part of the available sample (also called “test set” method). This “hold-out” is used to estimate the OOS performance, which is then compared with the IS performance. If they are congruent, the investor has grounds to “reject” the hypothesis that the backtest is overfit. The main advantage of this procedure is its simplicity. However, this approach is unsatisfactory for multiple reasons.

First, if the data is publicly available, it is quite likely that the researcher has used the “hold-out” as part of the IS. Second, even if no “hold-out” data was used, any seasoned researcher knows well how financial variables performed over the OOS interval, and that information will be used in the strategy design, consciously or not (see Schorfheide and Wolpin [20]).

Third, hold-out is clearly inadequate for small samples. The IS will be too short to fit, and the OOS too short to conclude anything with sufficient confidence. Weiss and Kulikowski [24] argue that hold-out should not be applied to an analysis with less than 1000 observations. For example, if a strategy trades on a weekly basis, hold-out could not be used on backtests of less than 20 years. In the same line, Van Belle and Kerr [23] point out the high variance of hold-out’s estimation errors. If we are unlucky, the chosen

---

<sup>1</sup>The Sharpe ratio, or reward-to-variability ratio, can be defined as the expected excess return per unit of risk. For details, see Bailey and Lopez de Prado [2].

“hold-out” may be the one that refutes a valid strategy or supports an invalid strategy. Different “hold-outs” are likely to lead to different conclusions.

Fourth, even if the researcher counts with a large sample, the OOS analysis will consume a large amount of the sample to be conclusive, which is detrimental to the strategy’s design (see Hawkins [10]). If the OOS is taken from the end of a time series, we are losing the most recent observations, which often are the most representative going forward. If the OOS is taken from the beginning of the time series, the testing will be done on the least representative portion of the data.

Fifth, as long as the researcher tries more than one strategy configuration, overfitting is always present (see Bailey et al. [1] for a proof). The hold-out method does not take into account the number of trials attempted before selecting a particular strategy configuration, and consequently hold-out cannot correctly assess a backtest’s representativeness. Hold-out leaves the investor guessing to what degree the backtest is overfit. The answer to the question “is this backtest overfit?” is not a true or false, but a non-null probability that depends on the number of trials involved (an input ignored by hold-out). In this paper we will show a way to compute this probability.

**Another approach.** An approach popular among practitioners consists in modeling the underlying financial variable, generate random scenarios and measure the performance of the investment strategy on those scenarios. This presents the advantage of generating a distribution of outcomes, rather than relying on a single OOS performance estimate, as the “hold-out” method does. The disadvantages are that the model that generates random series of the underlying variable may also be overfit, may not contain all relevant statistical features, and it has to be customized to every variable (large development costs). Some retail trading platforms offer backtesting procedures based on this approach, such as random generation of tick data by fractal interpolation.

Several procedures have been proposed to determine whether an econometric model is overfit, see White [25], Romano et al. [19], Harvey et al. [8] for a discussion in the context of Econometric models. Essentially these methods propose a way to adjust the  $p$ -values of estimated regression coefficients to account for the multiplicity of trials. These are valuable approaches when the trading rule relies on an econometric specification. That is not generally the case, as discussed in Bailey et al. [1].

**Our proposal.** A generic, model-free and non-parametric approach to backtest overfitting would be useful. The main innovation of this paper is to propose a general framework that adapts recent advances in experimental mathematics, machine learning and decision theory to the very particular problem of assessing the representativeness of a backtest. This is not an easy problem, as evidenced by the scarcity of academic papers addressing a dilemma that most investors face. This gap in the literature is surprising, considering practitioners’ reliance on backtests. One advantage of our solution is that it only requires time series of backtested performance. We avoid the credibility problem (of preserving a truly OOS test-set) by not requiring a fixed “hold-out,” and swapping all IS and OOS sets. Our approach is generic in the sense of not requiring knowledge of the trading rule or forecasting equation. The output is a bootstrapped distribution of OOS Sharpe ratios, as well as a measure of the representativeness of the backtest that we call probability of backtest overfitting (PBO). Although in our examples we always choose the Sharpe ratio to evaluate performance, our methodology can be applied to any other performance measure.

## 1.1 STRUCTURE OF THE PAPER

The rest of the study is organized as follows: Section 2 sets the foundations of our framework. Section 3 defines the PBO and some other useful statistics that can be derived from this approach. Section 4 discusses some of the features of our framework, and how it relates to other machine learning methods. Section 5 lists some of the limitations of this method. Section 6 presents several test cases to illustrate how the PBO compares to different scenarios. Section 7 assesses the accuracy of our method using two alternative approaches: Monte Carlo and Extreme Value Theory. Section 8 discusses a practical application. Section 9 summarizes our conclusions. The mathematical appendices prove the propositions presented throughout the paper.

## 2 THE FRAMEWORK

### 2.1 DEFINITION OF OVERFITTING IN THE CONTEXT OF STRATEGY SELECTION

We ask the question: What is the probability that an “optimal” strategy is overfit? Intuitively, for overfitting to occur, the strategy configuration that delivers maximum performance IS must systematically underperform

the rest of configurations OOS. The reason for this under performance is that the IS “optimal” strategy is too closely tied to the training set, to the point that optimizing becomes detrimental.

Consider a probability space  $(\mathcal{T}, \mathcal{F}, Prob)$  where  $\mathcal{T}$  represents a sample space of pairs of IS and OOS samples. We aim at estimating the probability of overfitting for the following *backtest strategy selection process*: select from  $N$  strategies labeled as  $(1, 2, \dots, N)$  the ‘best’ one using backtest according to a given performance measure, say, the Sharpe ratio. Fixing a performance measure, we will use random vectors  $\mathbf{R} = (R_1, R_2, \dots, R_N)$  and  $\bar{\mathbf{R}} = (\bar{R}_1, \bar{R}_2, \dots, \bar{R}_N)$  on  $(\mathcal{T}, \mathcal{F}, Prob)$  to represent the IS and OOS performance of the  $N$  strategies, respectively. For a given sample  $c \in \mathcal{T}$ , that is a concrete pair of IS and OOS samples, we will use  $\mathbf{R}^c$  and  $\bar{\mathbf{R}}^c$  to signify the performances of the  $N$  strategies on the IS and OOS pair given by  $c$ . For most applications  $\mathcal{T}$  will be finite and one can choose to use the power set  $\mathcal{T}$  as  $\mathcal{F}$ . Moreover, often it makes sense in this case to assume that the  $Prob$  is uniform on elements in  $\mathcal{T}$ . However, we do not make specific assumptions at this stage of general discussion so as to leave flexibilities in particular applications.

To discuss the probability of overfitting, the key is to compare the ranking of the selected strategies IS and OOS. Therefore we consider the ranking space  $\Omega$  consists of the  $N!$  permutations of  $(1, 2, \dots, N)$  indicating the ranking of the  $N$  strategies. Then we use random vectors  $r, \bar{r}$  to represent the ranking of the components of  $\mathbf{R}, \bar{\mathbf{R}}$ , respectively. For example, if  $N = 3$  and the performance measure is Sharpe ratio, for a particular sample  $c \in \mathcal{T}$ ,  $\mathbf{R}^c = (0.5, 1.1, 0.7)$  and  $\bar{\mathbf{R}}^c = (0.6, 0.7, 1.3)$  then we have  $r^c = (1, 3, 2)$  and  $\bar{r}^c = (1, 2, 3)$ . Thus, both  $r$  and  $\bar{r}$  are random vectors mapping  $(\mathcal{T}, \mathcal{F}, Prob)$  to  $\Omega$ .

Now, we define backtest overfitting, in the context of investment strategy selection alluded to above. We will need to use the following subset of  $\Omega$ :  $\Omega_n^* = \{f \in \Omega \mid f_n = N\}$ .

**Definition 2.1.** (Backtest Overfitting) *We say that the backtest strategy selection process overfits if a strategy with optimal performance IS has an expected ranking below the median OOS. By the Bayesian formula and using the notation above that is*

$$\sum_{n=1}^N E[\bar{r}_n \mid r \in \Omega_n^*] Prob[r \in \Omega_n^*] \leq N/2. \quad (2.1)$$

**Definition 2.2.** (Probability of Backtest Overfitting) *A strategy with optimal performance IS is not necessarily optimal OOS. Moreover, there is a*

non-null probability that this strategy with optimal performance IS ranks below the median OOS. This is what we define as the probability of backtest overfit (PBO). More precisely

$$PBO = \sum_{n=1}^N \text{Prob}[\overline{r}_n < N/2 \mid r \in \Omega_n^*] \text{Prob}[r \in \Omega_n^*] \quad (2.2)$$

In other words, we say that a strategy selection process overfits if the expected performance of the strategies selected IS is less than the median performance rank OOS of all strategies. In that situation, the strategy selection process becomes in fact detrimental. Note that in this context IS corresponds to the subset of observations used to select the optimal strategy among the  $N$  alternatives. With IS we do not mean the period on which the investment model underlying the strategy was estimated (e.g., the period on which crossing moving averages are computed, or a forecasting regression model is estimated). Consequently, in the above definition we refer to overfitting in relation to the strategy selection process, not a strategy's model calibration (e.g., in the context of regressions). That is the reason we were able to define overfitting without knowledge of the strategy's underlying models, i.e. in a model-free and non-parametric manner.

Estimating the probability of backtest overfitting in a particular application relies on schemes of selecting samples of IS and OOS pairs. Next, we describe a procedure to estimate the overfitting probability that we name *combinatorially symmetric cross-validation (CSCV)* for convenience of reference.

## 2.2 THE CSCV PROCEDURE

Suppose that a researcher is developing an investment strategy. She considers a family of system specifications and parametric values to be backtested, in an attempt to uncover the most profitable incarnation of that idea. For example, in a trend following moving average strategy, the researcher could try alternative sample lengths on which the moving averages are computed, entry thresholds, exit thresholds, stop losses, holding periods, sampling frequencies, etc. As a result, the researcher ends up running a number  $N$  of alternative model configurations (or trials), out of which one is chosen according to some performance evaluation criterion, such as the Sharpe ratio.

**Algorithm 2.3** (CSCV). We proceed as follows.

**First**, we form a matrix  $\mathbf{M}$  by collecting the performance series from the  $N$  trials. In particular, each column  $n = 1, \dots, N$  represents a vector of

profits and losses over  $t = 1, \dots, T$  observations associated with a particular model configuration tried by the researcher.  $\mathbf{M}$  is therefore a real-valued matrix of order  $(T \times N)$ . The only conditions we impose are that: i)  $\mathbf{M}$  is a true matrix, i.e. with the same number of rows for each column, where observations are synchronous for every row across the  $N$  trials, and ii) the performance evaluation metric used to choose the “optimal” strategy can be estimated on subsamples of each column. For example, if that metric was the Sharpe ratio, we would expect that the IID Normal distribution assumption can be held on various slices of the reported performance. If different model configurations trade with different frequencies, observations are aggregated to match a common index  $t = 1, \dots, T$ .

**Second**, we partition  $\mathbf{M}$  across rows, into an even number  $S$  of disjoint submatrices of equal dimensions. Each of these submatrices  $\mathbf{M}_s$ , with  $s = 1, \dots, S$ , is of order  $(T/S \times N)$ .

**Third**, we form all combinations  $C_S$  of  $\mathbf{M}_s$ , taken in groups of size  $S/2$ . This gives a total number of combinations

$$\binom{S}{S/2} = \binom{S-1}{S/2-1} \frac{S}{S/2} = \dots = \prod_{i=0}^{S/2-1} \frac{S-i}{S/2-i} \quad (2.3)$$

For instance, if  $S = 16$ , we will be forming 12,780 combinations. Each combination  $c \in C_S$  is composed of  $S/2$  submatrices  $\mathbf{M}_s$ .

**Fourth**, for each combination  $c \in C_S$ , we:

- a) Form the *training set*  $J$ , by joining the  $S/2$  submatrices  $\mathbf{M}_s$  that constitute  $c$  in their original order.  $J$  is a matrix of order  $(T/S)(S/2) \times N = T/2 \times N$ .
- b) Form the *testing set*  $\bar{J}$ , as the complement of  $J$  in  $M$ . In other words,  $\bar{J}$  is the  $T/2 \times N$  matrix formed by all rows of  $M$  that are not part of  $J$  also in their original order. (The order in forming  $J$  and  $\bar{J}$  does not matter for some performance measures such as the Sharpe ratio but does matter for others e.g. return maximum drawdown ratio).
- c) Form a vector  $\mathbf{R}^c$  of performance statistics of order  $N$ , where the  $n$ th component  $R_n^c$  of  $\mathbf{R}^c$  reports the performance associated with the  $n$ th column of  $J$  (the testing set). As before rank of the components of  $\mathbf{R}^c$  is denoted by  $r^c$  the IS ranking of the  $N$  strategies.
- d) Repeat c) with  $J$  replaced by  $\bar{J}$  (the test set) to derive  $\bar{\mathbf{R}}^c$  and  $\bar{r}^c$  the OOS performance statistics and rank of the  $N$  strategies, respectively.



- e) Determine the element  $n^*$  such that  $r_{n^*}^c \in \Omega_{n^*}^*$ . In other words,  $n^*$  is the best performing strategy IS.
- f) Define the relative rank of  $\bar{r}_{n^*}^c$  by  $\bar{\omega}_c := \bar{r}_{n^*}^c / (N + 1) \in (0, 1)$ . This is the relative rank of the OOS performance associated with the strategy chosen IS. If the strategy optimization procedure is not overfitting, we should observe that  $\bar{r}_{n^*}^c$  systematically outperforms OOS, just as  $r_{n^*}^c$  outperformed IS.
- g) We define the logit  $\lambda_c = \ln \bar{\omega}_c / (1 - \bar{\omega}_c)$ . High logit values imply a consistency between IS and OOS performances, which indicates a low level of backtest overfitting.

**Fifth**, compute the distribution of ranks OOS by collecting all the  $\lambda_c$ , for  $c \in C_S$ . Define the relative frequency at which  $\lambda$  occurred across all  $C_S$  by

$$f(\lambda) = \sum_{c \in C_S} \frac{\chi_{\{\lambda\}}(\lambda_c)}{\#(C_S)}, \quad (2.4)$$

where  $\chi$  is the characterization function and  $\#(C_S)$  signifies the number of elements in  $C_S$ . Then  $\int_{-\infty}^{\infty} f(\lambda) d\lambda = 1$ . This concludes the procedure.

Figure 1 schematically represents how the combinations in  $C_S$  are used to produce training and testing sets, where  $S = 4$ . It shows the six combinations of four subsamples A, B, C, D, grouped in two subsets of size two. The first subset is the training set (or in-sample). This is used to determine the optimal model configuration. The second subset is the testing set (or out-of-sample), on which the in-sample optimal model configuration is tested. Running the  $N$  model configurations over each of these combinations allows us to derive a relative ranking, expressed as a logit. The outcome is a distribution of logits, one per combination. Note that each training subset combination is re-used as a testing subset and vice-versa.

### 3 OVERFIT STATISTICS

The framework introduced in Section 2 allows us to characterize the reliability of a strategy's backtest in terms of four complementary analyses:

1. Probability of Backtest Overfitting (PBO): The probability that the model configuration selected as optimal IS will underperform the median of the  $N$  model configurations OOS.

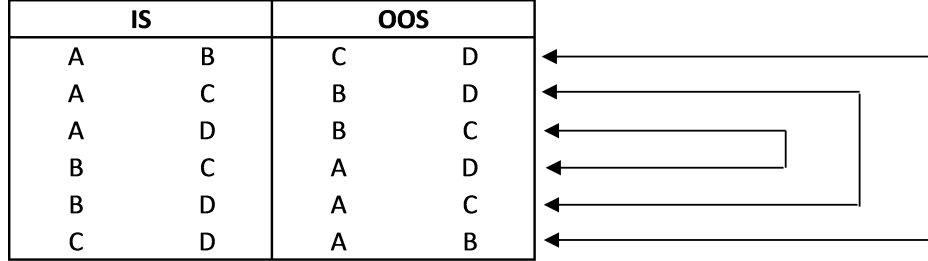


Figure 1: Generating the  $C_S$  symmetric combination

2. Performance degradation: It determines to what extent greater performance IS leads to lower performance OOS, an occurrence associated with the memory effects discussed in Bailey et al. [1].
3. Probability of loss: The probability that the model selected as optimal IS will deliver a loss OOS.
4. Stochastic dominance: This analysis determines whether the procedure used to select a strategy IS is preferable to randomly choosing one model configuration among the  $N$  alternatives.

### 3.1 PROBABILITY OF OVERFITTING (PBO)

PBO defined in Section 2.1 now can be estimated as  $\phi = \int_{-\infty}^0 f(\lambda) d\lambda$ . This represents the rate at which optimal IS strategies underperform the median of the OOS trials. The analogue of  $\bar{r}$  in medical research is the placebo given to a portion of patients in the test set. If the backtest is truly helpful, the optimal strategy selected IS should outperform most of the  $N$  trials OOS. That is the case when  $\lambda_c > 0$ . There are three relevant scenarios associated with  $\phi$ :

- $\phi \approx 0$ : A low proportion of combinations  $C$  exhibited  $\lambda_c < 0$ . Thus, the optimal IS strategy outperformed the median of trials in most of the testing sets. There is no significant overfitting, because choosing the optimal strategy IS indeed helped improve performance OOS.
- $\phi = 1/2$ : In half of the combinations  $C$  it occurred that  $\lambda_c < 0$ . The optimal IS strategy underperformed in as many trials as it outperformed. The expected performance of the optimal strategy identified

by the backtest equals the median performance from the trials. Backtests are overfit to the point that the strategy selection procedure does not add value.

- $\phi \gg 1/2$ : In more than half of the combinations  $C$  it occurred that  $\lambda_c < 0$ . Thus, the optimal IS strategy underperformed the median of trials in more than half of the testing sets. The degree of overfitting is so prevalent that choosing the optimal strategy leads to worse expected performance than picking a strategy at random from the trials.

### 3.2 PERFORMANCE DEGRADATION AND PROBABILITY OF LOSS

Section 2.2 introduced the procedure to compute, among other results, the pair  $(R_{n^*}, \overline{R_{n^*}})$  for each combination  $c \in C_S$ . Note that while we know that  $R_{n^*}$  is the maximum among the components of  $\mathbf{R}$ ,  $\overline{R_{n^*}}$  is not necessarily the maximum among the components of  $\overline{\mathbf{R}}$ . Because we are trying every combination of  $\mathbf{M}_s$  taken in groups of size  $S/2$ , there is no reason to expect the distribution of  $\mathbf{R}$  to dominate over  $\overline{\mathbf{R}}$ . The implication is that, generally,  $\overline{R_{n^*}} < \max\{\mathbf{R}\} \approx \max\{\overline{\mathbf{R}}\} = R_{n^*}$ . For a regression  $\overline{R_{n^*}}^c = \alpha + \beta R_{n^*}^c + \varepsilon^c$ , the  $\beta$  will be negative in most practical cases, due to compensation effects described in Bailey et al. [1]. An intuitive explanation for this negative slope is that overfit backtests minimize future performance: The model is so fit to the past, that it is rendered unfit for the future. And the more overfit a backtest is, the more memory is accumulated against its future performance.

It is interesting to plot the pairs  $(R_{n^*}, \overline{R_{n^*}})$  to visualize how strong is the performance degradation, and obtain a more realistic range of attainable performance OOS (see Figure 8). A particularly useful statistic is the proportion of combinations with negative performance,  $Prob[\overline{R_{n^*}}^c < 0]$ . Note that, even if  $\phi \approx 0$ ,  $Prob[\overline{R_{n^*}}^c < 0]$  could be high, in which case the strategy's performance OOS is poor for reasons other than overfitting.

Figure 2 provides a graphical representation of i) Out-Of-Sample Performance Degradation, ii) Out-Of-Sample Probability of Loss, and iii) Probability of Overfitting (PBO).

The upper plot of Figure 2 shows that pairs of (SR IS, SR OOS) for the optimal model configurations selected for each subset  $c \in C_S$ , which corresponding to the performance degradation associated with the backtest of an investment strategy. We can once again appreciate the negative relationship between greater SR IS and SR OOS, which indicates that at some point seeking the optimal performance becomes detrimental. Whereas 100% of the SR IS are positive, about 78% of the SR OOS are negative. Also, SR

## Figures

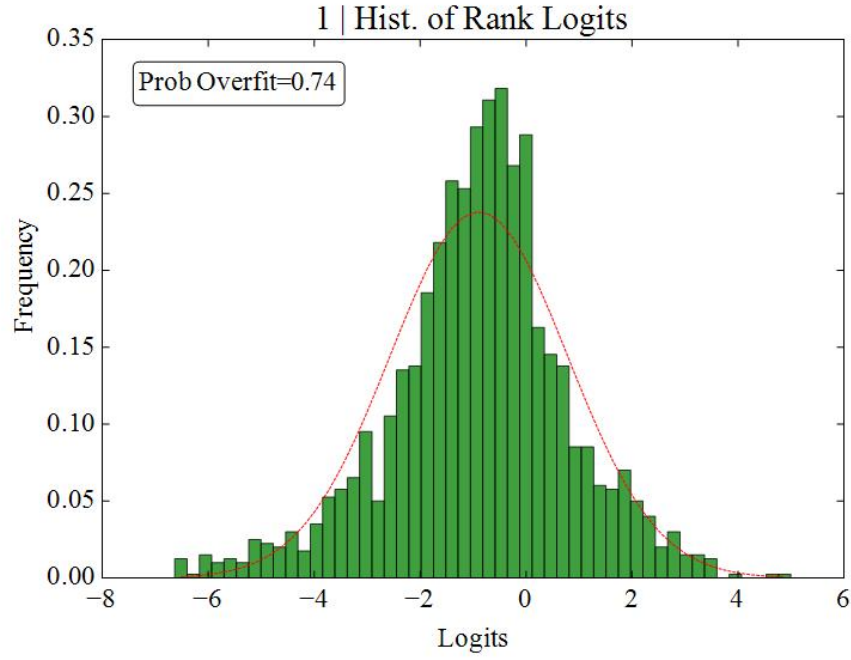
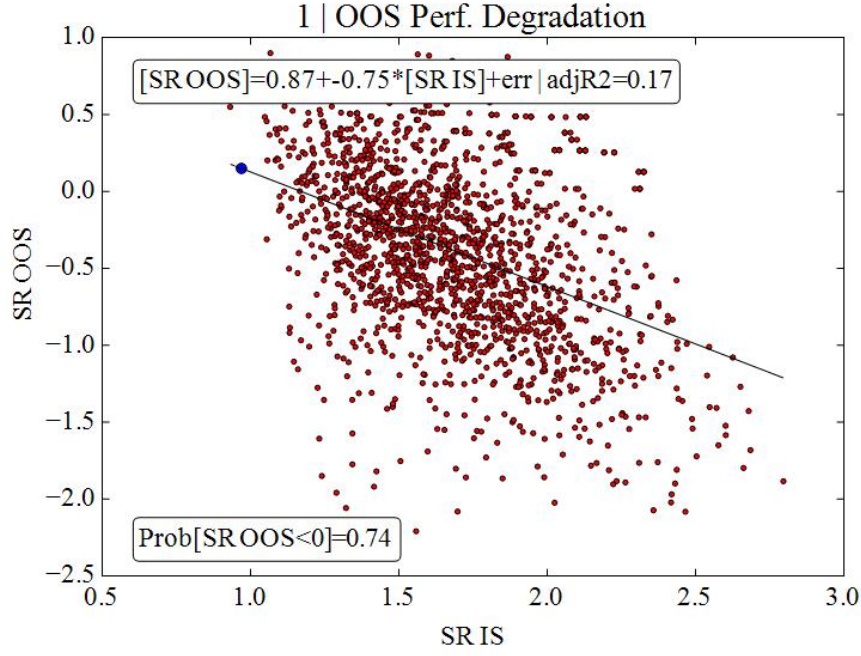


Figure 2: Performance degradation and distribution of logits  
 Note that, even if  $\phi \approx 0$ ,  $Prob[\overline{R}_n^c < 0]$  could be high, in which case the strategy's performance OOS is poor for reasons other than overfitting.

IS range between 1 and 3, indicating that backtests with high Sharpe ratios tell us nothing regarding the representativeness of that result.

We cannot hope escaping the risk of overfitting by exceeding some SR IS threshold. On the contrary, it appears that the higher the SR IS, the lower the SR OOS. In this example we are evaluating performance using the Sharpe ratio, however we stress that our procedure is generic and can be applied to any performance evaluation metric  $\mathbf{R}$  (Sortino ratio, Jensen's Alpha, Probabilistic Sharpe Ratio, etc.). This also allows us to compute the proportion of combinations with negative performance,  $Prob[\overline{R_{n^*}}^c < 0]$ , which corresponds to analysis ii).

The lower plot of Figure 2 shows the distribution of logits for the same strategy, with a PBO of 74%. It displays the distribution of logits, which allows us to compute the probability of backtest overfitting (PBO). This represents the rate at which optimal IS strategies underperform the median of the OOS trials.

Figure 3 plots the performance degradation and distribution of logits of a real investment strategy. Unlike in the previous example, the OOS probability of loss is very small (about 3%), and the proportion of selected (IS) model configurations that performed OOS below the median of overall model configurations was only 4%.

The upper plot of Figure 3 plots the performance degradation associated with the backtest of a real investment strategy. The regression line that goes through the pairs of (SR IS, SR OOS) is much less steep, and only 3% of the SR OOS are negative. The lower plot of Figure 3 shows the distribution of logits, with a PBO of 0.04%. According to this analysis, it is unlikely that this backtest is overfit. The chances that this strategy performs well OOS are much greater than in the previous example.

### 3.3 STOCHASTIC DOMINANCE

A further application of the results derived in Section 2.2 is to determine whether the distribution of  $\overline{R_{n^*}}$  across all  $c \in C_S$  stochastically dominates over the distribution of all  $\mathbf{R}$ . Should that not be the case, it would present strong evidence that strategy selection optimization does not provide consistently better OOS results than a random strategy selection. One reason that makes the concept of stochastic dominance useful is that it allows us to rank gambles or lotteries without having to make strong assumptions regarding an individual's utility function. See Hadar and Russell [6] for an introduction.

In the context of our framework, first-order stochastic dominance oc-

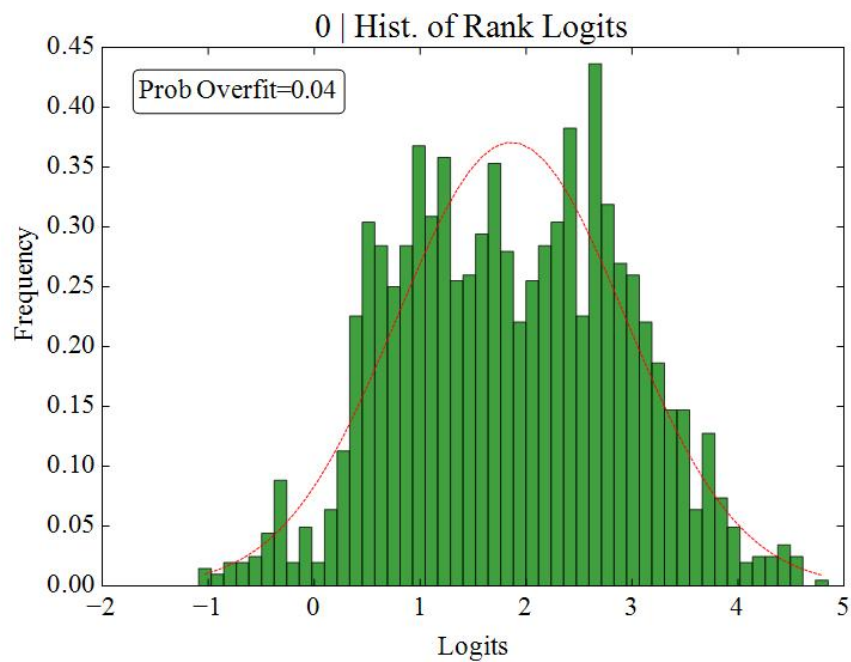
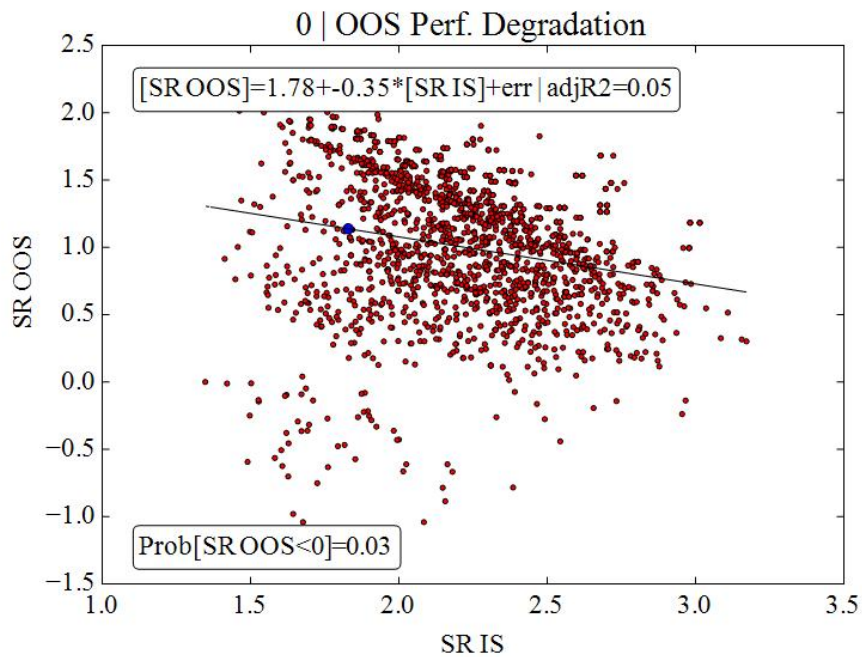


Figure 3: Performance degradation and distribution of logits for a real investment strategy

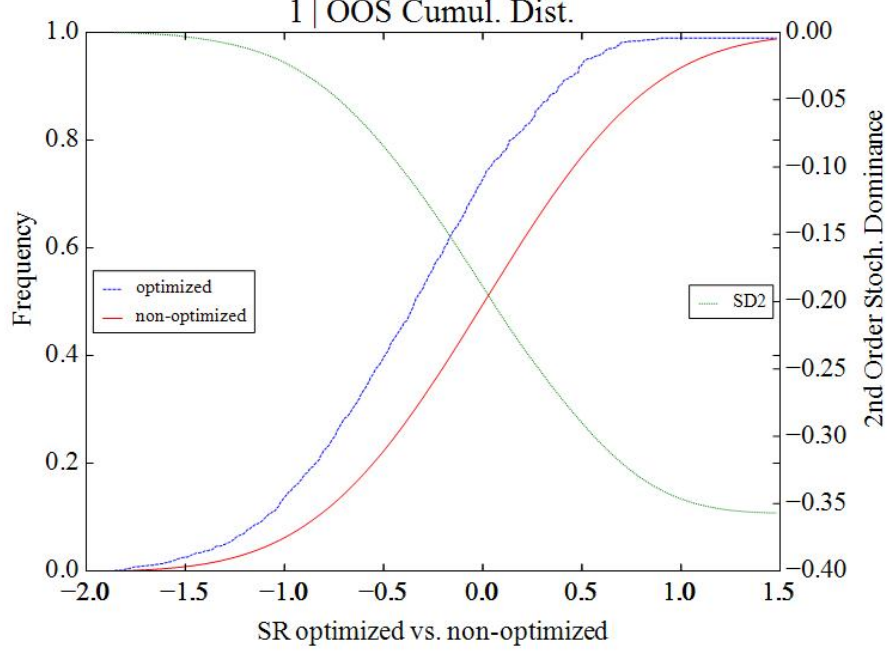


Figure 4: Stochastic dominance (example 1)

curs if  $Prob[\overline{R}_{n^*} \geq x] \geq Prob[Mean(\overline{\mathbf{R}}) \geq x]$  for all  $x$ , and for some  $x$ ,  $Prob[\overline{R}_{n^*} \geq x] > Prob[Mean(\overline{\mathbf{R}}) \geq x]$ . It can be verified visually by checking that the cumulative distribution function of  $\overline{R}_{n^*}$  is not above the cumulative distribution function of  $\overline{\mathbf{R}}$  for all possible outcomes, and at least for one outcome the former is strictly below the latter. Under such circumstances, the decision maker would prefer the criterion used to produce  $\overline{R}_{n^*}$  over a random sampling of  $\overline{\mathbf{R}}$ , as long as her utility function is weakly increasing.

A less demanding criterion is second-order stochastic dominance. This requires that  $SD2[x] = \int_{-\infty}^x (Prob[Mean(\overline{\mathbf{R}}) \leq x] - Prob[\overline{R}_{n^*} \leq x])dx \geq 0$  for all  $x$ , and that  $SD2[x] > 0$  at some  $x$ . When that is the case, the decision maker would prefer the criterion used to produce  $\overline{R}_{n^*}$  over a random sampling of  $\overline{\mathbf{R}}$ , as long as she is risk averse and her utility function is weakly increasing.

Figure 4 complements the analysis presented in Figure 2, with analysis of stochastic dominance. Stochastic dominance allows us to rank gambles or lotteries without having to make strong assumptions regarding an indi-

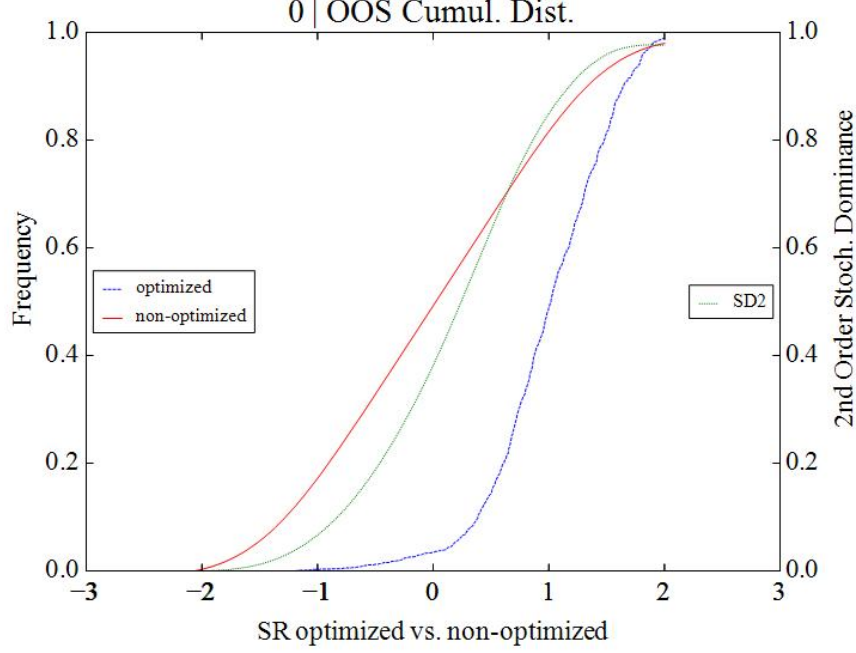


Figure 5: Stochastic dominance (example 2)

vidual's utility function.

Figure 4 also provides an example of the cumulative distribution function of  $\bar{R}_{n^*}$  across all  $c \in C_S$  (red line) and  $\bar{\mathbf{R}}$  (blue line), as well as the second order stochastic dominance ( $SD2[x]$ , green line) for every OOS SR. In this example, the distribution of OOS SR of optimized (IS) model configurations does not dominate (in first order) the distribution of OOS SR of overall model configurations. This can be appreciated in the fact that for every level of OOS SR the proportion of optimized model configurations is greater than the proportion of non-optimized, thus the probabilistic mass of the former is shifted to the left of the non-optimized.  $SD2$  plots the second order stochastic dominance, which indicates that the distribution of optimized model configurations does not dominate the non-optimized according to this less demanding criterion. It has been computed on the same backtest used for Figure 2. Consistent with that result, the overall distribution of OOS performance dominates the OOS performance of the optimal strategy selection procedure, a clear sign of overfitting.



Figure 5 provides a counter-example, based on the same real investment strategy used in Figure 6. It indicates that the strategy selection procedure used in this backtest actually added value, as the distribution of OOS performance for the selected strategies clearly dominates the overall distribution of OOS performance. First-order stochastic dominance is a sufficient condition for second-order stochastic dominance, and the plot of  $SD2[x]$  is consistent with that fact.

## 4 FEATURES OF THE CSCV SAMPLING METHOD

Our testing method combines multiple discoveries in the fields of machine learning (combinatorial optimization, jackknife, cross-validation) and decision theory (logistic function, stochastic dominance). Standard cross-validation methods include *k-fold cross-validation* (K-FCV) and *leave-one-out cross-validation* (LOOCV).

Now, K-FCV randomly divides the sample of size  $T$  into  $k$  subsamples of size  $T/k$ . Then it sequentially tests on each of the  $k$  samples the model trained on the  $T - T/k$  sample. Although a very valid approach in many situations, we believe that our procedure is more adequate than K-FCV in the context of strategy selection. In particular, we would like to compute the Sharpe ratio (or any other performance measure) on each of the  $k$  testing sets of size  $T/k$ . This means that  $k$  must be sufficiently small, so that the Sharpe ratio estimate is reliable (see Bailey and López de Prado [2] for a discussion of Sharpe ratio confidence bands). But if  $k$  is small, K-FCV will essentially reduce to a “hold-out” method, which we have argued is inaccurate.

Also, LOOCV is a K-FCV where  $k = T$ . We are not aware of any reliable performance metric computed on a single OOS observation.

The *combinatorially symmetric cross-validation* (CSCV) method we have proposed in Section 2.2 differs from K-FCV and LOOCV. The key idea is to generate  $\binom{S}{S/2}$  testing sets of size  $T/2$  by recombining the  $S$  slices of the overall sample of size  $T$ . This procedure presents a number of advantages. First, CSCV ensures that the training and testing sets are of equal size, thus providing comparable accuracy to the IS and OOS Sharpe ratios (or any other performance metric that is susceptible to sample size).

This is important, because making the testing set smaller than the training set (as hold-out does) would mean that we are evaluating with less accuracy OOS than the one used to choose the optimal strategy. Second, CSCV is symmetric, in the sense that all training sets are re-used as testing sets and vice versa. In this way, the decline in performance can only result from

overfitting, not arbitrary discrepancies between the training and testing sets.

Third, CSCV respects the time-dependence and other seasonalities present in the data, because it does not require a random allocation of the observations to the  $S$  subsamples. We avoid that requirement by recombining the  $S$  subsamples into the  $\binom{S}{S/2}$  testing sets. Fourth, CSCV derives a non-random distribution of logits, in the sense that each logit is deterministically derived from one item in the set of combinations  $C_S$ . Similarly to jackknife resampling, running CSCV twice on the same inputs generates identical results. Therefore, for each analysis, CSCV will provide a single result,  $\phi$ , which can be independently replicated and verified by another user. Fifth, the dispersion of the distribution of logits conveys relevant information regarding the robustness of the strategy selection procedure. A robust strategy selection leads to a consistent OOS performance rankings, which translate into similar logits.

Sixth, our procedure to estimate PBO is model-free, in the sense that it does not require the researcher to specify a forecasting model or the definitions of forecasting errors. It is also non-parametric, as we are not making distributional assumptions on PBO. This is accomplished by using the concept of logit,  $\lambda_c$ . A logit is the logarithm of odds. In our problem, the odds are represented by relative ranks (i.e., the odds that the optimal strategy chosen IS happens to underperform OOS). The logit function presents the advantage of being the inverse of the sigmoidal logistic distribution, which resembles the cumulative Normal distribution.

As a consequence, if  $\bar{\omega}_c$  are distributed close to uniform (the case when the backtest appears to be informationless), the distribution of the logits will approximate the standard Normal. This is important, because it gives us a baseline of what to expect in the threshold case where the backtest does not seem to provide any insight into the OOS performance. If good backtesting results are conducive to good OOS performance, the distribution of logits will be centered in a significantly positive value, and its left tail will marginally cover the region of negative logit values, making  $\phi \approx 0$ .

A key parameter of our procedure is the value of  $S$ . This regulates the number of submatrices  $M_s$  that will be generated, each of order  $(T/S \times N)$ , and also the number of logit values that will be computed,  $\binom{S}{S/2}$ .  $S$  must be large enough so that the number of combinations suffices to draw inference. If  $S$  is too small, the left tail of the distribution of logits will be underrepresented. On the other hand, if we believe that the performance series is time-dependent and incorporates seasonal effects,  $S$  cannot be too large, or the relevant time structure may be shuttered across the partitions.

For example, if the backtest includes more than six years of data,  $S = 24$

generates partitions spanning over a quarter each, which would preserve daily, weekly and monthly effects, while producing a distribution of 2,704,156 logits. In contrast, if we are interested in quarterly effects, we have two choices: i) Work with  $S = 12$  partitions, which will give us 924 logits, and/or ii) double  $T$ , so that  $S$  does not need to be reduced.

Another key parameter is the number of trials (i.e., the number of columns in  $M_s$ ). Hold-out's disregard for the number of trials attempted was the reason we concluded it was an inappropriate method to assess a backtest's representativeness (see Bailey et al. [1] for a proof).  $N$  must be large enough to provide sufficient granularity to the values of the relative rank,  $\bar{\omega}_c$ . If  $N$  is too small,  $\bar{\omega}_c$  will only adopt a very discrete number of values, which will translate in a very discrete number of logits, making  $f(\lambda)$  too discontinuous, adding estimation error to the evaluation of  $\phi$ . For example, if the investor is sensitive to values of  $\phi < 1/10$ , it is clear that the range of values that the logits can adopt must be greater than 10, and so  $N \gg 10$  is required. Other considerations regarding  $N$  will be discussed in the following Section.

Finally, PBO is evaluated by comparing combinations of  $T/2$  observations with their complementary. But the backtest counts with  $T$  observations, rather than only  $T/2$ . Therefore,  $T$  should be chosen to be double of the number of observations used by the investor to choose a model configuration or determine a forecasting specification.

## 5 LIMITATIONS AND MISUSE

This procedure was designed to evaluate PBO under minimal assumptions and input requirements. In doing so, we have attempted to provide a very general (in fact, model-free and non-parametric) procedure against which IS backtests can be benchmarked. The reader should understand that model-specific methods may be more accurate in certain instances, by sacrificing generality and comparability across models. For example, if a forecasting equation was used to generate the trials, it would be possible to develop a framework that evaluates PBO particular to that forecasting equation.

Because many investment strategies lack a forecasting equation, this is not always an option. We believe that a general procedure is useful in the context of deciding across multiple investment options, as exemplified by the success of benchmark methods such as the Sharpe ratio. But like in the case of the Sharpe ratio, it is important to discuss the limitations of our approach.

First, the researcher must provide as many profits-and-losses series ( $N$ ) as truly tested, and test as many strategy configurations as reasonable and feasible. Hiding trials will lead to an underestimation of the overfit, because each logit will be evaluated under a biased relative rank  $\overline{\omega_c}$ . It would be the equivalent to removing subjects from the trials of a new drug, once we have verified that the drug was not effective on them. Likewise, adding trials that are doomed to fail in order to make one particular model configuration succeed biases the result. If a model configuration is obviously flawed, it should have never been tried in the first place. A case in point are guided searches, where an optimization algorithm uses information from prior iterations to decide what direction should be followed next. In this case, the columns of matrix  $M$  should be the final outcome of each guided search (i.e., after it has converged to a solution), and not the intermediate steps.<sup>2</sup> This procedure aims at evaluating how reliable a backtest selection process is when choosing among feasible strategy configurations. As a rule of thumb, the researcher should backtest as many theoretically reasonable strategy configurations as possible.

Second, this procedure does nothing to evaluate the correctness of a backtest. If the backtest is flawed due to incorrect assumptions, such as transaction costs or using data not available at the moment of making a decision, our approach will be making an assessment based on flawed information.

Third, this procedure only takes into account structural breaks as long as they are present in the dataset of length  $T$ . If the structural break occurs outside the boundaries of the available dataset, the strategy may be overfit to a particular data regime, which our PBO has failed to account for because the entire set belongs to the same regime. This invites the more general warning that the dataset used for any backtest is expected to be representative of future states of the modeled financial variable.

Fourth, although a high PBO indicates overfitting in the group of  $N$  tested strategies, skillful strategies can still exist in these  $N$  strategies. For example, it is entirely possible that all the  $N$  strategies are having high but similar Sharpe ratios. Since none of the strategies is clearly better than the rest, PBO will be high. Here overfitting is among many skillful strategies.

Fifth, we must warn the reader against applying CSCV to guide the search for an optimal strategy. That would constitute a gross misuse of our method. As Strathern [22] eloquently put it, "when a measure becomes a

---

<sup>2</sup>We thank David Aronson and Timothy Masters (Baruch College) for asking for this clarification.

target, it ceases to be a good measure.” Any counter-overfitting technique used to select an optimal strategy will result in overfitting. For example, CSCV can be employed to evaluate the quality of a strategy selection process, but PBO should not be the objective function on which such selection relies.

## 6 TEST CASES

We will compare how PBO responds to several test cases: Full, high and low overfit. These cases are created by setting a matrix  $\mathbf{M}$  with  $N - 1$  trials of length  $T$  and null Sharpe ratio. If we add to that matrix  $\mathbf{M}$  a trial with Sharpe ratio zero, the strategy selection procedure will deliver a full overfit, because selecting the strategy with highest Sharpe ratio IS in this context cannot improve performance OOS over the median (zero). If we add to  $\mathbf{M}$  a trial with Sharpe ratio 1, selecting the strategy with the highest Sharpe ratio IS may still improve performance OOS over the median (zero), giving us the high overfit case. If we add to  $\mathbf{M}$  a trial with Sharpe ratio 2, selecting the strategy with the highest Sharpe ratio IS will likely improve performance OOS over the median (zero), giving us the low overfit case.

### 6.1 TEST CASE 1: FULL OVERFIT

We create a matrix  $\mathbf{M}$  where  $T = 1000, N = 100$  as follows.

1. For each trial,  $n = 1, \dots, 100$ :
  - a. Form a vector of  $T$  random draws from a standard Normal distribution.
  - b. Re-scale and re-centre the vector so that its Sharpe ratio is 0.
2. Combine the  $N$  vectors into a matrix  $\mathbf{M}$ .

If we choose IS the trial with highest Sharpe ratio, this cannot be expected to outperform OOS. If our procedure is accurate, we should obtain that  $\phi \rightarrow 1$ , indicating that in virtually all cases the “optimal” strategy IS happened to underperform OOS the median of trials. In effect, our simulations show that CSCV correctly determined that backtests are almost certain to overfit in this situation.

Increasing the sample size to  $T = 2000$  has no effect; we still obtain  $\phi \rightarrow 1$ . Increasing the number of trials to  $N = 200$  still produces  $\phi \rightarrow 1$ .

## 6.2 TEST CASE 2: HIGH OVERFIT

We create a matrix  $\mathbf{M}$  where  $T = 1000, N = 100$  as follows.

1. For each trial,  $n = 1, \dots, 100$ :
  - a. Form a vector of  $T$  random draws from a standard Normal distribution.
  - b. Re-scale and re-centre that vector so that its Sharpe ratio is 0.
2. Re-scale and re-centre the  $N$ th vector so that its Sharpe ratio is 1.
3. Combine the  $N$  vectors into a matrix  $\mathbf{M}$ .

Because one of the trials is expected to succeed OOS, PBO is not 1. At the same time, a random sample with a Sharpe ratio of 0 over  $T$  observations is likely to produce IS a Sharpe ratio above 1 over  $T/2$  observations (see Bailey et al. [1]). Accordingly, it is very likely that the strategy selection procedure will choose one trial with a high Sharpe ratio IS, only to underperform OOS the median of trials. The conclusion is that PBO will still be high in this scenario. Our Monte Carlo simulations confirm that intuition, by computing overfitting probabilities between 0.7 and 0.8.

Increasing the number of trials to  $N = 200$  slightly increases  $\phi$ , which now ranges between 0.75 and 0.85. The reason is, as more trials are added, the risk of overfitting increases, thus the importance that the researcher reports all the possibilities truly tested.

Increasing the sample size to  $T = 2000$  reduces PBO to values between 0.4 and 0.5, because the larger the number of available observations, the more informative is the performance IS. These empirical findings are consistent with Proposition 1 and Theorem 1, discussed in Bailey et al. [1].

## 6.3 TEST CASE 3: LOW OVERFIT

We create a matrix  $\mathbf{M}$  where  $T = 1000, N = 100$  as follows.

1. For each trial,  $n = 1, \dots, 100$ :
  - a. Form a vector of  $T$  random draws from a standard Normal distribution.
  - b. Re-scale and re-centre that vector so that its Sharpe ratio is 0.
2. Re-scale and re-centre the  $N$ th vector so that its Sharpe ratio is 2.

3. Combine the  $N$  vectors into a matrix  $\mathbf{M}$ .

Given that one of the trials performs significantly better than the rest over  $T$  observations, we expect a greater probability of it being selected IS over  $T/2$  observations. Still, there is a non-null probability that the strategy selection procedure fails to do so, because it is still possible for a subsample of that trial to underperform IS the subsample of one of the other trials, hence leading to overfitting. Accordingly, our simulations estimate overfitting probabilities ranging between 0.1 and 0.2.

Increasing the number of trials to  $N = 200$  slightly increases  $\phi$ , which now ranges between 0.2 and 0.3. Increasing the sample size to  $T = 2000$  reduces PBO to values between 0 and 0.04, just as we argued in Bailey et al. [1].

The reader may draw a parallel between these results and her knowledge of overfitting in regression models. As is documented in most statistics textbooks, the probability of overfitting a regression model increases as the number of degrees of freedom decreases. Our  $\mathbf{M}$  matrix has its regression analogue in the  $\mathbf{X}$  matrix of explanatory (or exogenous) variables, shaped in  $T$  rows (time observations) and  $N$  factors (columns), used to fit some other explained variable. In a regression model of this type, the number of degrees of freedom is  $T - N$ . The larger  $T$  is, the more degrees of freedom, and the lower the risk of overfitting. Conversely, the larger  $N$  is, the fewer degrees of freedom and the higher the risk of overfitting. In conclusion, the PBO model behaves as we would have expected in the familiar case of regression models.

## 7 ACCURACY OF THE TEST

In the previous Section, as we increased the Sharpe ratio of the  $N$ th trial above the other trials in matrix  $\mathbf{M}$ , we observed a decrease in PBO. Decreasing  $N$  and increasing  $T$  had a similar effect. Thus, overfitting estimates in the test cases above seem reasonable in an ordinal sense. The question remains, does the actual estimate correspond to the probability of the selected strategy to underperform the median of trials OOS? We evaluate the accuracy of our CSCV procedure to determine PBO in two different ways: Via Monte Carlo (MC) simulations, and applying Extreme Value Theory (EVT).

## 7.1 ACCURACY AS VIEWED BY MONTE CARLO SIMULATION

In order to determine the MC accuracy of our PBO estimate, we have generated 1,000 matrices  $\mathbf{M}$  (experiments) for various test cases of order  $(T \times N) = (1000 \times 100)$ , and computed the proportion of experiments that yielded an OOS performance below the median. If our CSCV procedure is accurate, the PBO that we estimated by resampling slices of a single matrix  $\mathbf{M}$  should be close to the probability derived from generating 1,000 matrices  $\mathbf{M}$  and computing the proportion of IS optima that underperformed the median OOS.

Snippet 1 lists a code written in Python that estimates PBO via Monte Carlo.

```
#!/usr/bin/env python
# On 20130704 by lopezdeprado@lbl.gov
#-----
def testAccuracy_MC(sr_base,sr_case):
    # Test the accuracy of CSCV against hold-out
    # It generates numTrials random samples and directly computes the
    # proportion where OOS performance was below the median.
    length,numTrials,numMC=1000,100,1000
    pathOutput='H:/Studies/Quant #23/paper/'
    #1) Determine mu,sigma
    mu_base,sigma_base=sr_base/(365.25*5/7.),1/(365.25*5/7.)**.5
    mu_case,sigma_case=sr_case/(365.25*5/7.),1/(365.25*5/7.)**.5
    hist,probOverfit=[],0
    #2) Generate trials
    for m in range(numMC):
        for i in range(1,numTrials):
            j=np.array([gauss(0,1) for j in range(length)])
            j*=sigma_base/np.std(j) # re-scale
            j+=mu_base-np.mean(j) # re-center
            j=np.reshape(j,(j.shape[0],1))
            if i==1:pnl=np.copy(j)
            else:pnl=np.append(pnl,j,axis=1)
    #3) Add test case
    j=np.array([gauss(0,1) for j in range(length)])
    j*=sigma_case/np.std(j) # re-scale
    j+=mu_case-np.mean(j) # re-center
    j=np.reshape(j,(j.shape[0],1))
```



```

pnl=np.append(pnl,j,axis=1)
#4) Run test
# Reference distribution
mu_is=[np.average(pnl[:length/2,i]) for i in range(pnl.shape[1])]
sigma_is=[np.std(pnl[:length/2,i]) for i in range(pnl.shape[1])]
mu_oos=[np.average(pnl[length/2:,i]) for i in range(pnl.shape[1])]
sigma_oos=[np.std(pnl[length/2:,i]) for i in range(pnl.shape[1])]
sr_is=[mu_is[i]/sigma_is[i] for i in range(len(mu_is))]
sr_oos=[mu_oos[i]/sigma_oos[i] for i in range(len(mu_oos))]
print m,sr_is.index(max(sr_is)),max(sr_is), \
      sr_oos[sr_is.index(max(sr_is))]
sr_oos_=sr_oos[sr_is.index(max(sr_is))]
hist.append(sr_oos_)
if sr_oos_<np.median(sr_oos):probOverfit+=1
probOverfit/=float(numMC)
print probOverfit
return

```

Snippet 1 - Python code for estimating PBO via Monte Carlo

## 7.2 ACCURACY BY EXTREME VALUE THEORY

Backtesting a number  $N$  of alternative strategy configurations and selecting the trial that exhibits maximum performance IS sets the background for applying Extreme Value Theory. From [1, Eq. (2.3)] we know that Sharpe ratio estimates asymptotically follow a Gaussian distribution. Proposition 1 in [1] discussed the distribution of the maximum of  $N$  independent random variables. As part of its proof, we applied the Fisher-Tippett-Gnedenko theorem to the Gaussian distribution, and concluded that the maximum performance IS among  $N$  alternative backtests can be approximated through a Gumbel distribution.

More formally, suppose a set of backtests where  $N = 100, T = 1000$ ,  $SR_n = 0$  for  $n = 1, \dots, N - 1$  and  $SR_N = \widetilde{SR} > 0$ . The sample length is divided in two sets of equal size  $T/2$ , IS and OOS. A strategy with  $SR_n = 0$  is selected when its IS SR is greater than the IS SR of the strategy with  $SR_N = \widetilde{SR} > 0$ . Because the Sharpe ratio has been globally constrained for the whole sample by re-scaling and re-centering, and IS has the same length as OOS,  $SR_{OS}^* \approx SR - SR_I S^*$ . By virtue of this global constraint, the following propositions can be used to estimate PBO:

- i. The median of all Sharpe ratios OOS is null,  $Me[SR_{OOS}] = 0$ .

- ii. Selecting a strategy with  $SR_n = 0$  leads to  $SR_{OOS}^* \approx \max_N < Me[SR_{OOS}]$  iff  $SR_{IS}^* > 0$ .
- iii. Selecting a strategy with  $SR_N = \widetilde{SR}$  leads to  $SR_{OOS}^* \approx \widetilde{SR} - SR_{IS}^*$ , where  $E[SR_{OOS}^*] > Me[SR_{OOS}]$  iff  $SR_{IS}^* \in (-\infty, 2\widetilde{SR})$  and  $E[SR_{OOS}^*] \leq Me[SR_{OOS}]$  iff  $SR_{IS}^* \in [2\widetilde{SR}, \infty)$ .
- iv. Computing  $SR_{IS}$  fully determines  $SR_{OOS}$  as a result of the global constraint. Thus,  $V[SR_{IS}] = V[SR] = (1 + SR^2/2)/T$ , and  $V[SR_{OOS}] = 0$ .

Our goal is to compute the probability that the strategy with maximum Sharpe ratio IS performs below the median of OOS Sharpe ratios. First, we need to calibrate the parameters of the Gumbel distribution associated with a set of Gaussian random variables. Suppose a random variable  $\max_N = \max(\{SR_n | n = 1, \dots, N-1\})$ , where  $SR_n$  is the Sharpe ratio estimated through a backtest for trial  $n$ . We know that the Gaussian distribution belongs to the Maximum Domain of Attraction of the Gumbel distribution, thus  $\max_N \sim \Lambda[\alpha, \beta]$ , where  $\alpha, \beta$  are the normalizing constants and  $\Lambda$  is the CDF of the Gumbel distribution. Next, we derive the values of these normalizing constants. It is known that the mean and standard deviation of a Gumbel distribution are

$$\begin{aligned} E[\max_N] &= \alpha + \gamma\beta \\ \sigma[\max_N] &= \frac{\beta\pi}{\sqrt{6}} \end{aligned} \tag{7.1}$$

where  $\gamma$  is the Euler-Mascheroni constant. Applying the method of moments, we can derive the following:

- Given an estimate of  $\sigma[\max_N]$ , we can estimate  $\hat{\beta} = \frac{\hat{\sigma}[\max_N]\sqrt{6}}{\pi}$ .
- Given an estimate of  $E[\max_N]$ , and the previously obtained  $\hat{\beta}$ , we can estimate  $\hat{\alpha} = \hat{E}[\max_N] - \gamma\hat{\beta}$ .

These parameters allow us to model the distribution of the maximum Sharpe ratio IS out of a set of  $N-1$  trials. PBO can then be directly computed as  $\phi = \phi_1 + \phi_2$ , with

$$\begin{aligned} \phi_1 &= \int_{-\infty}^{2\widetilde{SR}} \mathcal{N}\left[SR, \frac{1+SR^2/2}{T}\right] (1 - \Lambda[\alpha(SR), \beta(SR)]) dSR \\ \phi_2 &= \int_{2\widetilde{SR}}^{\infty} \mathcal{N}\left[SR, \frac{1+SR^2/2}{T}\right] dSR \end{aligned} \tag{7.2}$$

where  $\beta(SR) = \sqrt{6}(1 + \frac{1}{2}SR^2)/(\pi T)$  and  $\alpha(SR) = \max(0, SR) - \gamma\beta(SR)$  are the estimate of  $\beta$  and  $\alpha$  based on (7.1). Probability  $\phi_1$  accounts for selecting IS a strategy with  $SR_n = 0$ , as a result of  $SR_{N,IS} < SR_{IS}^*$ . As we argued earlier, in this situation  $SR_{OOS}^* - \max_N < Me[SR_{OOS}] = 0$  iif  $SR_{IS}^* > 0$ , hence the  $\max(0, SR)$  used to evaluate the Gumbel distribution. The integral has an upper boundary  $2\widetilde{SR}$  because beyond that point all trials lead to  $SR_{OOS}^* < Me[SR_{OOS}]$ , including the  $N$ th trial. That probability is accounted for by  $\phi_2$ , which has a lower boundary of integration in  $2\widetilde{SR}$ . Snippet 2 implements the numerical integration of Eq. (7.2). As we will see in Section 7.3, the EVT estimates of PBO derived from Eq. (7.2) are in agreement with the MC estimates.

```
#!/usr/bin/env python
# On 20130704 by lopezdeprado@lbl.gov
#-----
def testAccuracy_EVT(sr_base, sr_case):
    # Test accuracy by numerical integration
    # It does the same as testAccuracy_MC, but through numerical integration ...
    # ... of the base and case distributions.
    #1) Parameters
    parts,length,freq,minX,trials=1e4,1000,365.25*5/7.,-10,100
    emc=0.57721566490153286 # Euler-Mascheroni constant
    #2) SR distributions
    dist_base=[sr_base,((freq+.5*sr_base**2)/length)**.5]
    dist_case=(sr_case,((freq+.5*sr_case**2)/length)**.5)
    #3) Fit Gumbel (method of moments)
    maxList=[]
    for x in range(int(parts)):
        max_=max([gauss(dist_base[0],dist_base[1]) for i in range(trials)])
        maxList.append(max_)
    dist_base[1]=np.std(maxList)*6**.5/math.pi
    dist_base[0]=np.mean(maxList)-emc*dist_base[1]
    #4) Integration
    prob1=0
    for x in np.linspace(minX*dist_case[1],2*dist_case[0]-sr_base,parts):
        f_x=ss.norm.pdf(x,dist_case[0],dist_case[1])
        F_y=1-ss.gumbel_r.cdf(x,dist_base[0],dist_base[1])
        prob1+=f_x*F_y
    prob1*=(2*dist_case[0]-sr_base-minX*dist_case[1])/parts
    prob2=1-ss.norm.cdf(2*dist_case[0]-sr_base,dist_case[0],dist_case[1])
```

```

print dist_base,dist_case
print prob1,prob2,prob1+prob2
return

```

Snippet 2 - Python code for computing PBO via EVT

### 7.3 AN EMPIRICAL STUDY OF ACCURACY

We are finally ready to evaluate the accuracy of CSCV’s PBO. To achieve that, we will compare CSCV’s PBO estimates against the two alternative benchmarks described in Sections 7.1 (MC) and 7.2 (EVT). Table 2 reports the results for a wide range of combinations of  $\widetilde{SR}$  (SR\_Case),  $T$  and  $N$ . Without loss of generality,  $SR_n = 0$  for  $n = 1, \dots, N - 1$ . We do not need to consider alternative values of  $SR_n$ , because the likelihood of selecting the wrong strategy is a function of  $\widetilde{SR} - SR_n$ , not the absolute level of  $SR_n$ .

The table in Figure 6 shows PBO estimates using three alternative methods: Combinatorially Symmetric Cross-Validation (CSCV), Monte Carlo (MC) and Extreme Value Theory (EVT).

Monte Carlo results were computed on 1,000 experiments. The proportion of IS optimal selections that underperformed OOS is reported in Prob\_MC. Column Prob\_EVT reports the corresponding PBO estimates, derived from Eq. (7.2). Because these latter results are derived from the actual distribution of the maximum SR, they are more accurate than the Monte Carlo estimates. In any case, EVT and MC results are very close, with a maximum absolute deviation of 4.2%.

We have computed CSCV’s PBO on 1,000 randomly generated matrices  $\mathbf{M}$  for every parameter combination  $(\widetilde{SR}, T, N)$ . Therefore, we have obtained 1,000 independent estimates of PBO for every parameter combination, with a mean and standard deviation reported in columns Mean\_CSCV and Std\_CSCV. This is not to be mistaken with the Monte Carlo result, which produced a single estimate of PBO out of 1,000 randomly generated matrices  $\mathbf{M}$ .

A comparison of the Mean\_CSCV probability with the EVT result gives us an average absolute error of 2.1%, with a standard deviation of 2.9%. The maximum absolute error is 9.9%. That occurred for the combination  $(\widetilde{SR}, T, N) = (3, 500, 500)$ , whereby CSCV gave a more conservative estimate (24.7% instead of 14.8%). There is only one case where CSCV underestimated PBO, with an absolute error of 0.1%. The median error is only 0.7%, with a 5%-tile of 0% and a 95%-tile of 8.51%.

SR_Case	T	N	Mean_CSCV	Std_CSCV	Prob_MC	Prob_EVT	CSCV-EVT
0	500	500	1.000	0.000	1.000	1.000	0.000
0	1000	500	1.000	0.000	1.000	1.000	0.000
0	2500	500	1.000	0.000	1.000	1.000	0.000
0	500	100	1.000	0.000	1.000	1.000	0.000
0	1000	100	1.000	0.000	1.000	1.000	0.000
0	2500	100	1.000	0.000	1.000	1.000	0.000
0	500	50	1.000	0.000	1.000	1.000	0.000
0	1000	50	1.000	0.000	1.000	1.000	0.000
0	2500	50	1.000	0.000	1.000	1.000	0.000
0	500	10	1.000	0.001	1.000	1.000	0.000
0	1000	10	1.000	0.000	1.000	1.000	0.000
0	2500	10	1.000	0.000	1.000	1.000	0.000
1	500	500	0.993	0.007	0.991	0.994	-0.001
1	1000	500	0.893	0.032	0.872	0.870	0.023
1	2500	500	0.561	0.022	0.487	0.476	0.086
1	500	100	0.929	0.023	0.924	0.926	0.003
1	1000	100	0.755	0.034	0.743	0.713	0.042
1	2500	100	0.371	0.034	0.296	0.288	0.083
1	500	50	0.870	0.031	0.878	0.859	0.011
1	1000	50	0.666	0.035	0.628	0.626	0.041
1	2500	50	0.288	0.047	0.199	0.220	0.068
1	500	10	0.618	0.054	0.650	0.608	0.009
1	1000	10	0.399	0.054	0.354	0.360	0.039
1	2500	10	0.123	0.048	0.093	0.086	0.036
2	500	500	0.679	0.037	0.614	0.601	0.079
2	1000	500	0.301	0.038	0.213	0.204	0.097
2	2500	500	0.011	0.011	0.000	0.002	0.009
2	500	100	0.488	0.035	0.413	0.405	0.084
2	1000	100	0.163	0.045	0.098	0.099	0.065
2	2500	100	0.004	0.006	0.002	0.001	0.003
2	500	50	0.393	0.040	0.300	0.312	0.081
2	1000	50	0.113	0.044	0.068	0.066	0.047
2	2500	50	0.002	0.004	0.000	0.000	0.002
2	500	10	0.186	0.054	0.146	0.137	0.049
2	1000	10	0.041	0.027	0.011	0.023	0.018
2	2500	10	0.000	0.001	0.000	0.000	0.000
3	500	500	0.247	0.043	0.174	0.148	0.099
3	1000	500	0.020	0.017	0.005	0.005	0.015
3	2500	500	0.000	0.000	0.000	0.000	0.000
3	500	100	0.124	0.042	0.075	0.068	0.056
3	1000	100	0.007	0.008	0.001	0.002	0.005
3	2500	100	0.000	0.000	0.000	0.000	0.000
3	500	50	0.088	0.037	0.048	0.045	0.043
3	1000	50	0.004	0.006	0.002	0.001	0.003
3	2500	50	0.000	0.000	0.000	0.000	0.000
3	500	10	0.028	0.022	0.010	0.015	0.013
3	1000	10	0.001	0.002	0.000	0.001	0.000
3	2500	10	0.000	0.000	0.000	0.000	0.000

Figure 6: CSCV's accuracy

In conclusion, CSCV provides accurate estimates of PBO, with relatively small errors on the conservative side.

## 8 A PRACTICAL APPLICATION

Bailey et al. [1] present an example of an investment strategy that attempts to profit from a seasonal effect. For the reader’s convenience, we reiterate here how the strategy works. Suppose that we would like to identify the optimal monthly trading rule, given four customary parameters: *Entry\_day*, *Holding\_period*, *Stop\_loss* and *Side*.

*Side* defines whether we will hold long or short positions on a monthly basis. *Entry\_day* determines the business day of the month when we enter a position. *Holding\_period* gives the number of days that the position is held. *Stop\_loss* determines the size of the loss as a multiple of the series’ volatility that triggers an exit for that month’s position. For example, we could explore all nodes that span the interval  $[1, \dots, 22]$  for *Entry\_day*, the interval  $[1, \dots, 20]$  for *Holding\_period*, the interval  $[0, \dots, 10]$  for *Stop\_loss*, and  $\{-1, 1\}$  for *Sign*. The parameters combinations involved form a four-dimensional mesh of 8,800 elements. The optimal parameter combination can be discovered by computing the performance derived by each node.

As discussed in the above cited paper, a time series of 1000 daily prices (about 4 years) was generated by drawing from a random walk. Parameters were optimized (*Entry\_day* = 11, *Holding\_period* = 4, *Stop\_loss* = -1 and *Side* = 1), resulting in an annualized Sharpe ratio is 1.27. Given the elevated Sharpe ratio, we could conclude that this strategy’s performance is significantly greater than zero for any confidence level. Indeed, the PSR-Stat is 2.83, which implies a less than 1% probability that the true Sharpe ratio is below 0 (see Bailey and López de Prado [2] for details). Figure 7 is an graphical illustration of this example.

We have estimated the PBO using our CSCV procedure, and obtained the results illustrated below. Figure 8 shows that approx. 53% of the SR OOS are negative, despite all SR IS being positive and ranging between 1 and 2.2. Figure 9 plots the distribution of logits, which implies that, despite the elevated SR IS, the PBO is as high as 55%. Consequently, Figure 10 shows that the distribution of optimized OOS SR does not dominate the overall distribution of OOS SR. This is consistent with the fact that the underlying series follows a random walk, thus the serial independence among observations makes any seasonal patterns coincidental. The CSCV framework has succeeded in recognizing that the backtest was overfit.

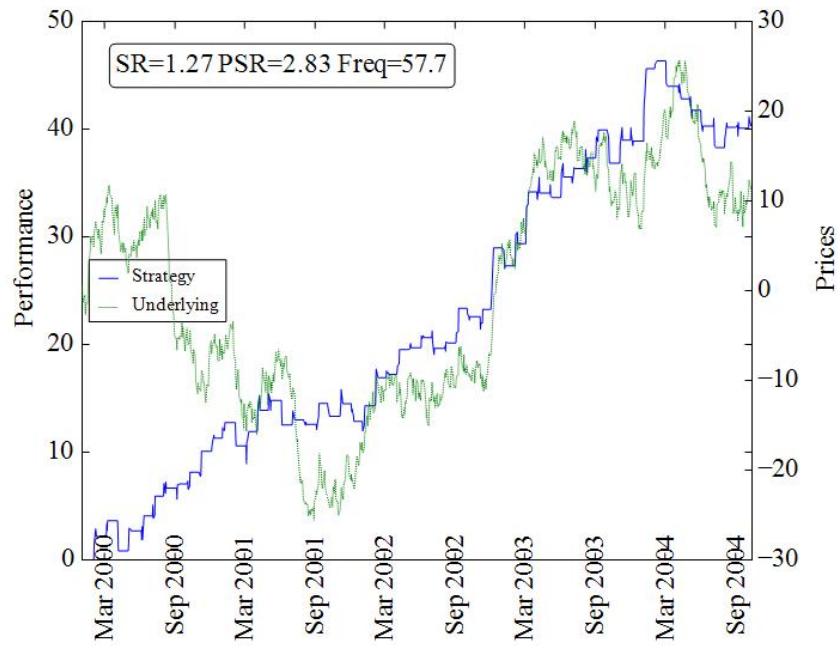


Figure 7: Backtested performance of a seasonal strategy (example 1)

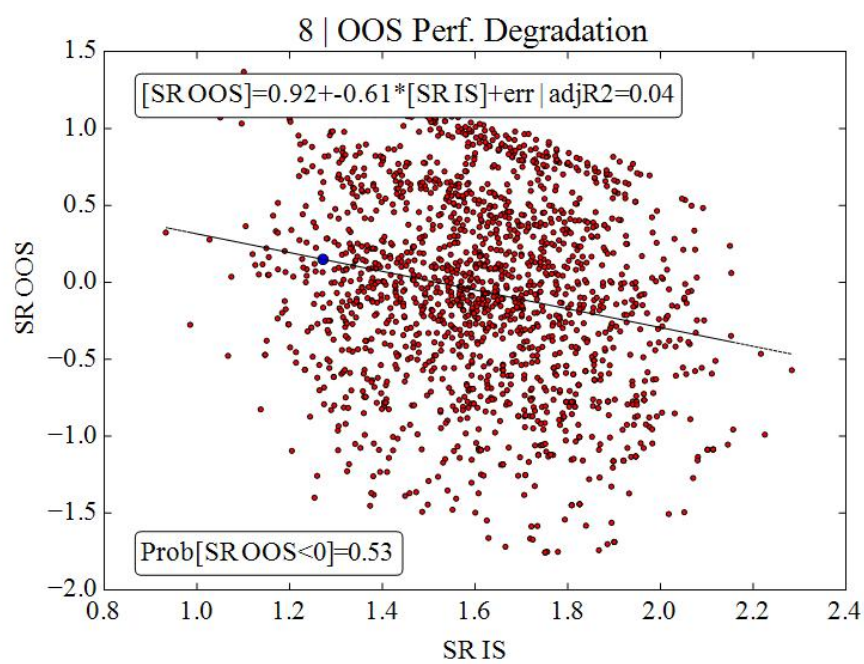


Figure 8: CSCV analysis of the backtest of a seasonal strategy (example 1): Performance degradation



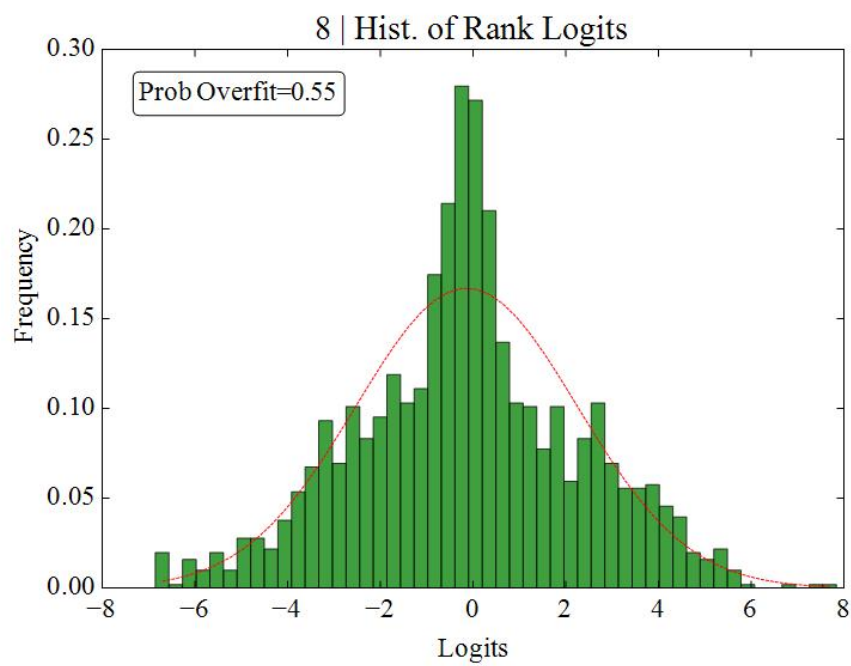


Figure 9: CSCV analysis of the backtest of a seasonal strategy (example 1): logit distrubution

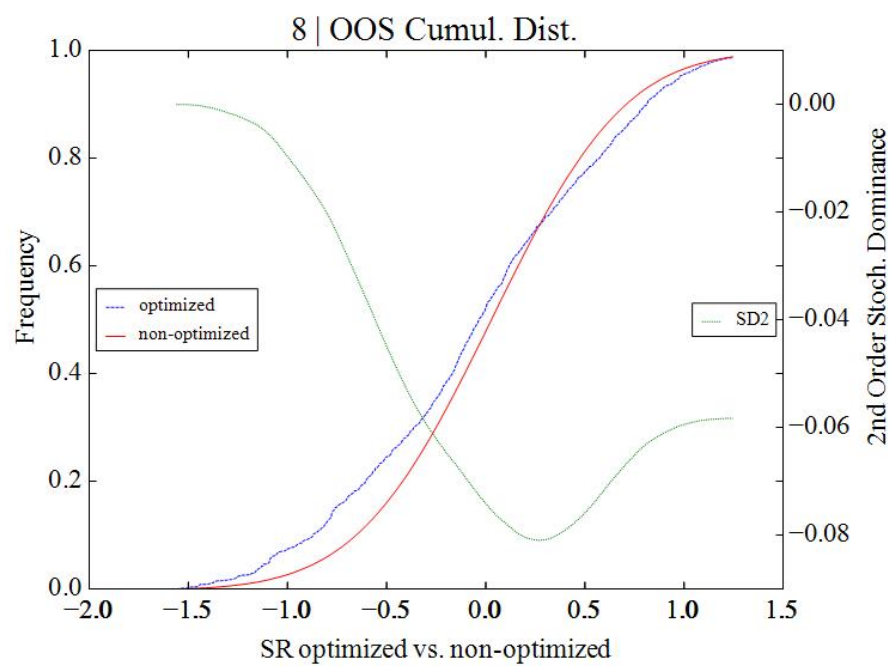


Figure 10: CSCV analysis of the backtest of a seasonal strategy (example 1): Absent of dominance

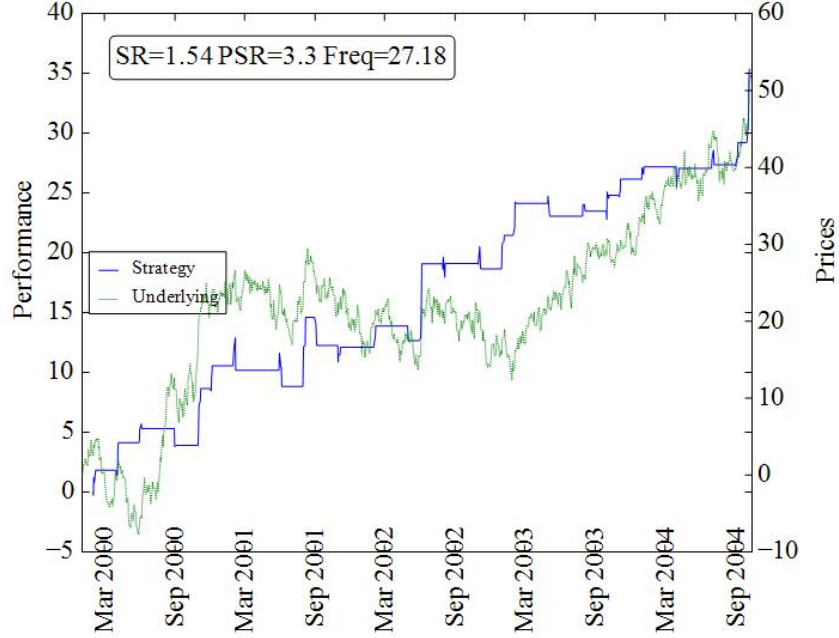


Figure 11: Backtested performance of a seasonal strategy (example 2)

Second, we generated a time series of 1000 daily prices (about 4 years), following a random walk. But unlike in the first case, we have shifted the returns of the first 5 random observations of each month to be centered at a quarter of a standard deviation. This generates a monthly seasonal effect, which the strategy selection procedure should discover. Figure 11 plots the random series, as well as the performance associated with the optimal parameter combination:  $\text{Entry\_day} = 1$ ,  $\text{Holding\_period} = 4$ ,  $\text{Stop\_loss} = -10$  and  $\text{Side} = 1$ . The annualized Sharpe ratio is 1.54 similar to the previous (overfit) case (1.54 vs. 1.3).

The next three graphs report the results of the CSCV analysis, which confirm the validity of this backtest in the sense that performance inflation from overfitting is minimal. Figure 12 shows that only 13% of the OOS SR to be negative. Because there is a real monthly effect in the data, the PBO for this second case should be substantially lower than the PBO of the first case. Figure 13 shows a distribution of logits with a PBO of only 13%. Figure 14 evidences that the distribution of OOS SR from IS optimal

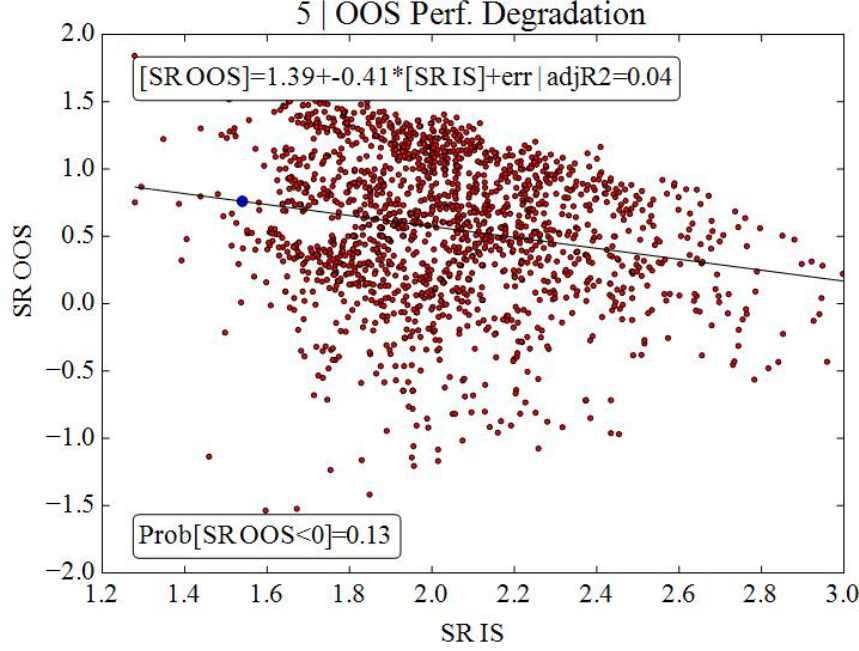


Figure 12: CSCV analysis of the backtest of a seasonal strategy (example 2): monthly effect

combinations clearly dominates the overall distribution of OOS SR. The CSCV analysis has correctly recognized the validity of this backtest, in the sense that performance inflation from overfitting is small.

In this practical application we have illustrated how simple is to produce overfit backtests when answering common investment questions, such as the presence of seasonal effects. We refer the reader to [1, Appendix 4] for the implementation of this experiment in Python language. Similar experiments can be designed to demonstrate overfitting in the context of other effects, such as trend-following, momentum, mean-reversion, event-driven effects, etc. Given the facility with which elevated Sharpe ratios can be manufactured IS, the reader would be well advised to remain critical of backtests and researchers that fail to report the PBO.

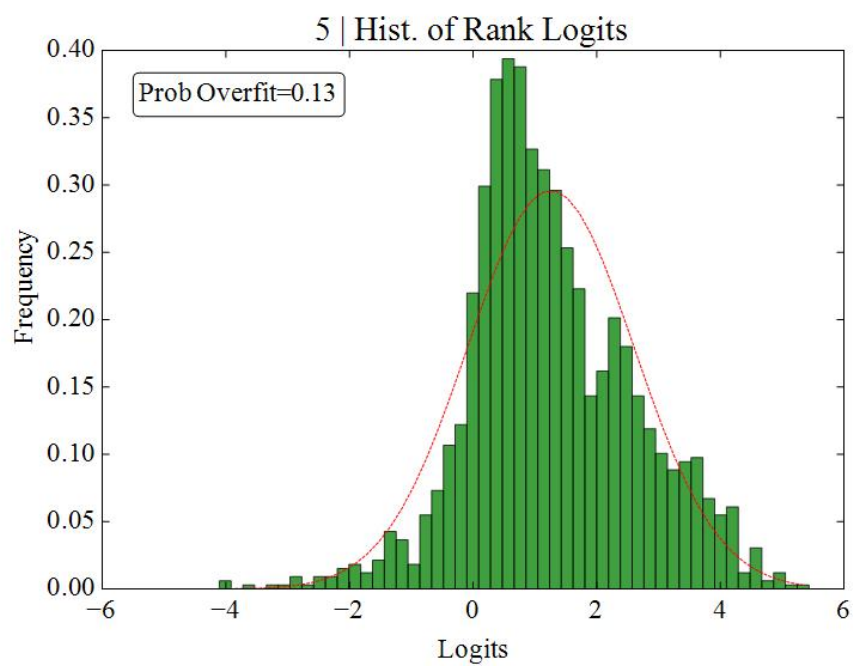


Figure 13: CSCV analysis of the backtest of a seasonal strategy (example 2): logit distribution

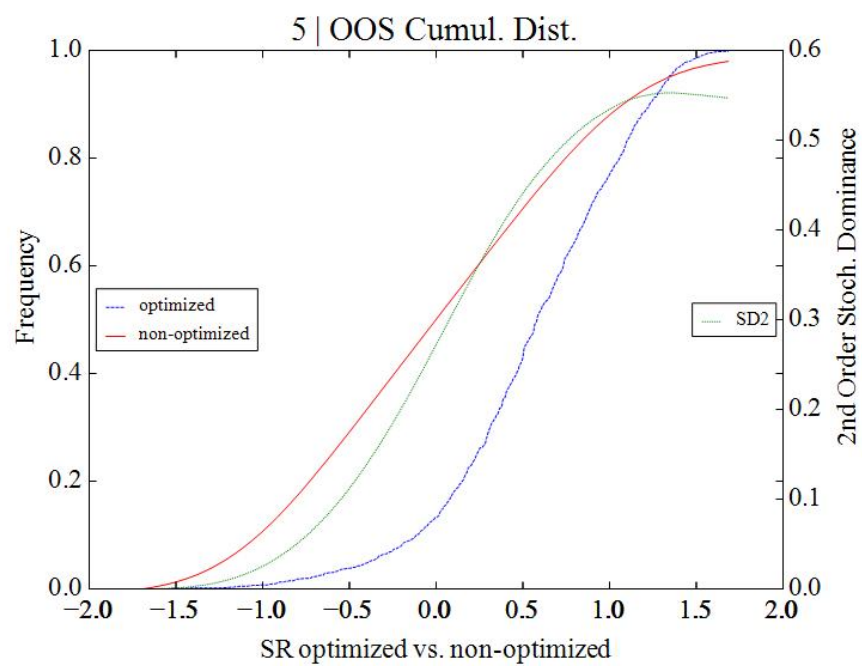


Figure 14: CSCV analysis of the backtest of a seasonal strategy (example 2): dominance

## 9 CONCLUSIONS

Bailey and López de Prado [2] developed methodologies to evaluate the probability that a Sharpe ratio is inflated (PSR), and the minimum track record length (MinTRL) required for a Sharpe ratio to be statistically significant. These statistics were developed to assess Sharpe ratios based on live performance (track record). We have not been able to find similar statistics or methodologies applicable to evaluate backtested Sharpe ratios. The representativeness of backtested performance estimates has been the subject of this study.

Standard statistical techniques designed to detect overfitting in the context of regression models are poorly equipped to assess backtest overfitting. Hold-outs in particular are unreliable and easy to manipulate. As long as a researcher tries more than one strategy configuration, overfitting is always present. However, hold-out methods do not take into account the number of trials involved in the strategy selection ( $N$ ), making it an inappropriate method to evaluate a backtest’s representativeness.

The procedure presented in this paper has specifically been designed to determine the probability of backtest overfitting (PBO). This is defined as the probability that the strategy with optimal performance IS delivers OOS a performance below the median performance of all trials attempted by the researcher. When that probability is high, optimizing IS has a detrimental effect in terms of OOS performance, because the backtest profits from specific features in the IS subset that are not present elsewhere. CSCV identifies such situation by generating a large number of combinations of IS subsets, and determining the proportion of those where overfitting has taken place. Unlike hold-out tests, CSCV takes into account the number of trials attempted ( $N$ ). Unlike in the case of hold-out tests, the decision as to whether overfitting occurs does not depend on an arbitrary division of the data into an IS and an OOS set, but on all possible combinations. In addition, all IS subsets are re-used as OOS, and all OOS subsets are re-used as IS.

We have assessed the accuracy of CSCV with regards to the estimation of PBO in two different ways, on a wide variety of test cases. Monte Carlo simulations show that CSCV applied on a single dataset provides similar results to computing PBO on a large number of independent samples. We have also computed directly PBO by deriving the Extreme Value distributions that model the performance of IS optimal strategies. Results indicate that CSCV provides accurate estimates of PBO, with relatively small errors on the conservative side.

In conclusion, we believe that CSCV is a new and powerful tool in the arsenal of investors and financial markets' researchers. At least we hope that this study raises the awareness concerning the futility of computing and reporting backtest results without controlling for its PBO and MinBTL.

## References

- [1] Bailey, D., J. Borwein, M. López de Prado and J. Zhu, "Pseudomathematics and financial charlatanism: The effects of backtest over fitting on out-of-sample performance," *Notices of the AMS*, to appear May 2014. Available at <http://ssrn.com/abstract=2308659>.
- [2] Bailey, D. and M. López de Prado, "The Sharpe Ratio Efficient Frontier," *Journal of Risk*, 15(2012), 3–44. Available at <http://ssrn.com/abstract=1821643>.
- [3] Doyle, J. and C. Chen, "The wandering weekday effect in major stock markets," *Journal of Banking and Finance*, 33 (2009), 1388–1399.
- [4] Embrechts, P., C. Klueppelberg and T. Mikosch, *Modelling Extremal Events*, Springer Verlag, New York, 2003.
- [5] Feynman, R., *The Character of Physical Law*, 1964, The MIT Press.
- [6] Hadar, J. and W. Russell, "Rules for Ordering Uncertain Prospects," *American Economic Review*, 59 (1969), 25–34.
- [7] Harris, L., *Trading and Exchanges: Market Microstructure for Practitioners*, Oxford University Press, 2003.
- [8] Harvey, C. and Y. Liu, "Backtesting", SSRN, working paper, 2013. Available at [http://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2345489](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2345489).
- [9] Harvey, C., Y. Liu and H. Zhu, "...and the Cross-Section of Expected Returns," SSRN, 2013. Available at [http://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2249314](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2249314).
- [10] Hawkins, D., "The problem of overfitting," *Journal of Chemical Information and Computer Science*, 44 (2004), 10–12.
- [11] Hirsch, Y., *Don't Sell Stocks on Monday*, Penguin Books, 1st Edition, 1987.



- [12] Leinweber, D. and K. Sisk, “Event Driven Trading and the ‘New News’,” *Journal of Portfolio Management*, 38(2011), 110–124.
- [13] Leontief, W., “Academic Economics”, *Science*, 9 Jul 1982, 104–107.
- [14] Lo, A., “The Statistics of Sharpe Ratios,” *Financial Analysts Journal*, 58 (2002), July/August.
- [15] López de Prado, M. and A. Peijan, “Measuring the Loss Potential of Hedge Fund Strategies,” *Journal of Alternative Investments*, 7 (2004), 7–31. Available at <http://ssrn.com/abstract=641702>.
- [16] López de Prado, M. and M. Foreman, “A Mixture of Gaussians Approach to Mathematical Portfolio Oversight: The EF3M Algorithm,” *Quantitative Finance*, to appear, 2014. Available at <http://ssrn.com/abstract=1931734>.
- [17] Mayer, J., K. Khairy and J. Howard (2010): “Drawing an Elephant with Four Complex Parameters,” *American Journal of Physics*, 78 (2010), 648–649.
- [18] Resnick, S., *Extreme Values, Regular Variation and Point Processes*, Springer, 1987.
- [19] Romano, J. and M. Wolf, “Stepwise multiple testing as formalized data snooping”, *Econometrica*, 73 (2005), 1273–1282.
- [20] Schorfheide, F. and K. Wolpin, “On the Use of Holdout Samples for Model Selection,” *American Economic Review*, 102 (2012), 477–481.
- [21] Stodden, V., Bailey, D., Borwein, J., LeVeque, R, Rider, W. and Stein, W., “Setting the default to reproducible: Reproducibility in computational and experimental mathematics,” February, 2013. Available at <http://www.davidhbailey.com/dhbpapers/icerm-report.pdf>.
- [22] Strathern, M., ”Improving Ratings: Audit in the British University System,” *European Review*, 5, (1997) pp. 305-308.
- [23] Van Belle, G. and K. Kerr, *Design and Analysis of Experiments in the Health Sciences*, John Wiley and Sons, 2012.
- [24] Weiss, S. and C. Kulikowski, *Computer Systems That Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning and Expert Systems*, Morgan Kaufman, 1st Edition, 1990.

- [25] White, H., “A Reality Check for Data Snooping,” *Econometrica*, 68 (2000), 1097–1126.
- [26] Wittgenstein, L., *Philosophical Investigations*, Blackwell Publishing, 1953. Section 201.