

**DISCUSSION ON THE PAPER: *ELICITABILITY AND BACKTESTING: PERSPECTIVES FOR BANKING REGULATION*
BY NATALIA NOLDE AND JOHANNA F. ZIEGEL**

BY MARIE KRATZ*

*ESSEC Business School, CREAR risk research center**

1. Context. Since the seminal paper of Artzner et al. ((3)), there has been a large literature on risk measures, their properties, their impact on practice and regulation. In the paper by Nolde and Ziegel, the authors review the specific properties of elicibility and backtestability of risk measures and their impacts in view of model validation and banking regulation.

The two popular risk measures used in financial institutions and regulation are Value-at-Risk (VaR) and Expected Shortfall (ES). VaR has been the dominant risk measure in banking regulation, although it is not a coherent risk measure for all distributions unlike ES (see (2), (18)). Backtesting those two risk measures depend on the type (point or distribution) forecast methods that have been used. For the VaR, a direct backtest is possible and the most popular procedure is based on the so-called violation process. In practice, it consists in replacing VaR at the $\alpha\%$ level by its estimates and checking that this process behaves like independent and identically distributed Bernoulli random variables with violation (success) probability close to $1 - \alpha$. If the proportion of VaR violations is not significantly different from $1 - \alpha$, the estimation/prediction method is concluded to be reasonable (see e.g. (8), (9) to test the independence assumption, (7) for a review on VaR backtesting procedures, and (13) for further references on the topic).

As a lesson of the last financial crisis, and notwithstanding the absence of an obvious direct backtest method for ES, it has been discussed then decided that a 10-day ES at the 97.5% level would replace the daily 99% VaR under Basel 3 (see (4)), ES being a more tail-sensitive measure of risk, hence more adequate for assessing extreme risks.

This decision caused some debate among all actors: academics, professionals of the banks and regulators, especially in terms of backtesting. First Gneiting ((15)) raised that there could be an issue with direct backtesting of ES estimates because ES is not elicitable, unlike VaR, an optimal point forecast

elicited by the weighted absolute error scoring function. Nevertheless, if not elicitable, ES has been proved to be conditionally elicitable ((13)) (which corresponds to the step procedure often done in practice) as well as jointly elicitable ((14)). Then, Acerbi & Székely (1) pointed out that elicibility (or lack of elicibility) is not relevant for backtesting of risk measures, but rather for comparing the forecast performance of different estimation methods. Currently, a consensus has emerged that the problem of comparing the forecasting performance of different estimation methods (which requires the property of 'elicibility') is distinct from the problem of checking whether a single series of *ex ante* estimates of a risk measure is consistent with a series of *ex post* realizations of P&L (*i.e.* 'backtestability'). Moreover there are reasonable approaches to backtest ES, some of them being developed recently to answer the concerns of banks; the simplest one, based on the violation process, suggests an implicit backtest for ES via the simultaneous backtest of four VaR (see (16) and (17)).

In this paper, the authors suggest to use coherent and elicitable risk measures in order to adopt a two-stage backtesting framework. The first one would correspond to the current practice. If passing the first stage, the second one would allow for a comparative backtest between standard and internal models.

2. The Study. The paper is developed into two parts. The first one suggests a theoretical framework to revisit the state of the art in terms of risk measures and the concepts of elicibility, identifiability, dominance (conditional and on average), before proposing comparative backtests that would complement (or supplement?) traditional backtests when using elicitable risk measures. The second part illustrates partly those tests on simulated and real data. Three elicitable risk measures are considered in all examples and applications: VaR, the pair (VaR, ES) and also the expectile.

I will come back on the main points, discussing their pros and cons in view of both the theory and its applications results.

- *Calibration tests*

In view of regulation, a one sided conditional super-calibration test is proposed component-wise (when considering a multivariate risk measure), to simplify the problem. For the joint risk measure (VaR, ES), choosing a component-wise test implies the equivalence between a conditional sub-calibration test for (VaR, ES) and testing that the conditional VaR, ES respectively, is at least as large as its optimal conditional prediction, which might appear a bit counter-intuitive. It illustrates how relations between

super (or sub) calibration and over (or under) estimation depend on the identification functions (and the component-wise condition). It would be interesting to see if using the concept of conditional elicitable risk measure ((13)) instead might lead to a more direct relation.

For the VaR risk measure, the conditional calibration test for $\mathbf{h}_t = 1$ is close to the standard backtest for VaR based on the violation process, and does not require, for one-step ahead forecasts, an asymptotic test (under adequate conditions) as provided; the exact Rosenblatt-Diebold-Davis test gives indeed the result, moreover with no condition. It is for multi-steps ahead forecasts that the proposed asymptotic test might be useful. Note that regulators do not ask for such forecasts, but for a qualitative evaluation of the risk evolution over the next three years.

In Table 2, we observe that H_0 is rejected more often with the two-sided test than with the one-sided (except for the (VaR,ES) at low threshold), as expected. General conditional calibration tests are more conservative than simple ones. For VaR and Expectile, general tests are able to better discriminate between models than simple ones when considering the first two thresholds; for the Basel reference threshold, such an observation is more contestable as the general test might reject H_0 when it should not, especially when using the expectile. It might be the problem of using this specific risk measure at such a high threshold. For the bidimensional risk measure (VaR, ES), the decision of rejecting H_0 does not really depend on whether the test is simple or general. In Table 3, when using real data (Nasdaq), we observe similar results as those obtained with simulation data: the choice of the distribution (likelihood function) seems to have a lower impact than the choice of the method at the second stage, which is not what we would like; again the expectile does not give reliable results with the general CCT; for (VaR,ES), the general CCT rejects all methods using a normal likelihood, even with the EVT method, but not with the simple test in this latter case.

In view of the obtained results, the main question remains *how to choose the predictable test functions* (using the available information at $t - 1$) for the conditionally calibrated test to be effective on finite samples. Right now, the test functions have been deduced from a parallel made with existing tests, on VaR (from Diebold-Davis) or on ES (from McNeil-Frey). It is always interesting to have a larger theoretical frame to cover specific cases, as suggested in this paper, but it still deserves further study to see if it can go beyond the existing, and especially improve current practice, in particular for regulation.

- *An alternative way for model validation*

If looking at the size of a traditional test is necessary, looking at its power is also crucial, especially for regulators, as emphasized in (17). But it is true that the alternative may be quite broad. Hence it is a very good idea to look for an alternative way to compute the power, introducing as here a comparative backtest on the prediction between two models with two possible null hypotheses, permuting the role of each model. The authors develop theoretically this backtest having in view a comparison between standard and internal model, but, when applying it numerically, come back to a comparison between classical trading book models calibrated with standard estimation methods, as done when computing the power for traditional test considering various possible models for the alternatives (see e.g. (10) or (11) when looking at probabilistic forecasts for the tail, and, in the case of point forecasts in the tail, (16) or (17) where the same type of experiments as in this paper has been performed). So the applications might have been chosen to better highlight the supplement of traditional backtests (defined with specific alternatives) with a comparative backtest.

Note also that considering the average $\Delta_n \bar{S}$ over time ($t = 1, \dots, n$) for comparative backtests reformulated according to Diebold-Mariano ((12)), might raise an issue in practice, if the score differences are not first order stationary! We will come back later on this issue.

- *Standard and Internal Models*

The standard model in a bank considers three types of risks, independently of each other: the credit risk, the market risk and the operational risk. When dealing with market risk, standard and internal models are based on the same principle, evaluating the VaR at 99% in Basel 2 and the ES 97.5% in Basel 3. It is essentially for the two other risks, credit and operational, that standard and internal models might differ, as the techniques used to measure them are usually more sophisticated with internal models than with standard ones.

Here the study is made on market risk (taking into account the daily assets), for which no real distinction is made between internal and standard models. So it appears a bit artificial, if not misleading, to use this terminology in the various figures where only estimation methods of the innovation distribution are compared. Other types of tests could answer the concern when computing the power, with no need for the risk measure to be elicitable or for the data to satisfy some conditions.

Nevertheless, I find it very interesting to propose a comparative backtest between standard and internal models, as addressed theoretically in this paper when using elicitable risk measures. It may be seen as a first step of an

alternative way for regulators and for the banks to check about the relevance of their internal model; it certainly deserves much more investigation to make it a useful tool.

- *Scoring*

In the numerical study, ranking between the methods via the average scores is given in Tables 1 (simulated data) and 3 (Nasdaq example). We observe that the ranking depends on the choice of the scoring function, with generally more stability for the VaR than for the two other risk measures. In Table 3, judging on the VaR and (VaR,ES), the scoring functions seem to be more sensitive to the estimation method than to the model.

Besides the problem of finding the 'right' scoring function, the main general question concerns their stability with time (stationarity). Considering the mean over t instead of comparing $E[S(R_t) - S(R_t^*)]$ for each t might raise some issues in practice and weaken a decision tool. The question has been passed over when speaking about the S-dominance on average and the hypothesis of the score differences as first-order stationary.

For the scoring to be useful, it needs to be a good predictor of the future ranking/scoring. Otherwise it could be really misleading and would even induce large errors relying on an unstable decisional tool. To illustrate this question, we could draw the parallel, for instance, with investment funds. Ranking has been made in terms of the funds performance (not in terms of risk measures), but it has been observed that the ranking is not at all a good predictor for the future ranking. The fund exhibiting the best performance one year will not be ranked number one the following year.

- *Choice of Risk Measures*

Besides the two 'regulatory' risk measures, VaR and ES, the authors consider another risk measure, the expectile. It is a coherent and elicitable risk measure (see (5), (19), (6)) but not comonotonically additive (see (13)). Expectiles, besides appearing less attractive for practitioners because of a less intuitive underlying concept than the concepts for VaR or ES, show clearly in this study some limits for their use in practice, specifically for regulation perspective. There is no explicit estimator of the expectile, whatever is the estimation method (Full Parametric, Filtered Historical or EVT). Moreover, to have a comparable magnitude of risk as VaR at 99% or ES at 97.5%, expectile requires to look at the 99.855% level, which might make statistical studies quite shaky, as observed on the results. In Table 2, if the expectile presents the best results at the lowest threshold (96.561%) when using a two-sided general test, the Student distribution is more often rejected than

the normal one at level 98.761% with this test, and, when looking at the 'Basel reference threshold', it is worse. The true model is rejected when using a two-sided general test for the FP and FHS methods but not rejected with the EVT one. We could argue that it makes sense at such high level, but we observe the reverse phenomenon when using the misspecified likelihood. So we cannot conclude to an over-sensitivity towards estimation methods rather than towards likelihoods at this high level. For the comparative backtest given in Table 1, the expectile behaves more or less as the other risk measures, exhibiting some instability depending on the choice of the scoring function. When considering the highest level, we observe some issue if one uses the FHS method, for the expectile with one given scoring function, and for the VaR with both scoring functions. In the example on real data, the expectile gives bad results, with no help to take any decision most of the time. Figures 1 to 4 provide a visible way to judge risk measures in terms of predictability. All figures confirm that expectile is not of great help for decision making. Various reasons might explain this, as for instance the choice of the test function, but it does not push in favor of the use of this risk measure in practice.

3. Conclusion. The authors have put forward the right theoretical questions and came up with interesting alternative backtests to compare models and methods. The numerical study performed on simulated data and a real data example, shows that there is still a long way to propose a practical tool. The authors themselves point that out.

For conditional calibration tests, the choice of the test functions is a bit tricky. Simulations and applications show that it still needs to be investigated to be effective in practice. I am in favor of rigorous but simple methods for practitioners, with not too many tuning parameters, in order to avoid most of the errors coming from using results relying on conditions not always easy to check out.

For the scoring functions as well, we observe some variability; further investigation is needed. But for the scoring to be useful, it has to be a good predictor of the future ranking/scoring. This is definitively the first problem to be looked at.

I would like to end this discussion with a general comment on how to push for the right incentive providing the right tools. The main issue is to optimize the capital, finding the right and fair amount for all actors: companies, shareholders, regulators, society at large. A company underestimating or overestimating her risks means that she does not manage them well, not

balancing the opposite goals of the various actors. It is true that the role of regulators is to check that a company did not underestimate her risks, not to put in danger her clients but also, as a consequence, society (that has been often asked to pay for the errors made by banks). I believe their incentive is as well to push companies to evaluate their risks accurately to foster good risk management. The new regulatory rules are going in this direction, encouraging companies to develop their internal models for better understanding their risks. The primary goal of regulation is clearly to protect consumers, but it is also important for regulation to favor best practices and encourage the development of a healthy industry. Hence I am not convinced that checking conditional super-calibration (or sub-calibration) is key to risk management. If a one-sided test might look at first glance reasonable from a regulatory point of view, a two-sided test is a useful test for companies, even if it has not been asked by regulators yet.

References.

- [1] ACERBI, C., SZÉKELY, B. (2014). Backtesting expected shortfall. *Risk magazine* **27**, 76–81.
- [2] ACERBI, C., TASCHE, D. (2002). On the coherence of expected shortfall. *Journal of Banking and Finance* **26**, 1487–1503.
- [3] ARTZNER, P., DELBAEN, F., EBER, J.-M. , HEATH, D. (1999). Coherent measures of risks. *Mathematical Finance* **9**, 203-228.
- [4] BASEL COMMITTEE ON BANKING SUPERVISION (2016). Minimum capital requirements for market risk. *Bank of International Settlements*.
- [5] BELLINI, F., KLAR, B., MÜLLER, A., AND ROSAZZA GIANIN, E. (2014). Generalized quantiles as risk measures. *Insurance: Mathematics and Economics* **54**, 41–48.
- [6] CHEN, J.M. (2013). Measuring Market Risk Under Basel II, 2.5, and III: VAR, Stressed VAR, and Expected Shortfall. Available at SSRN: <http://ssrn.com/abstract=2252463>
- [7] CAMPBELL, S.D. (2006). A review of Backtesting and Backtesting Procedures. *Journal of Risk* **9(2)**, 1–17.
- [8] CHRISTOFFERSEN, P. (2003). *Elements of Financial Risk Management*. Academic Press.
- [9] CHRISTOFFERSEN, P., AND PELLETIER, D. (2004). Backtesting Value-at-Risk: A Duration-Based Approach. *J. Financial Econometrics* **2(1)**, 84–108.
- [10] COSTANZINO, N., AND CURRAN, M. (2015). Backtesting general spectral risk measures with application to expected shortfall. *J. Risk Model Validation* **9**, 21–33.
- [11] COSTANZINO, N., AND CURRAN, M. (2016). A simple traffic light approach to backtesting expected shortfall. Working paper.
- [12] DIEBOLD, F.X., AND MARIANO, R.S. (1995). Comparing predictive accuracy. *Journal of Business and Economic Statistics* **13**, 253–263.
- [13] EMMER, S., KRATZ, M., AND TASCHE, D. (2015). What is the best risk measure in practice? *Journal of Risk* **18**, 31–60.
- [14] FISSLER, T., AND ZIEGEL, J.F. (2016). Higher order elicibility and Osband’s principle. *Annals of Statistics* **44**, 1680–1707.
- [15] GNEITING, T. (2011). Making and evaluating point forecasts. *Journal of the American Statistical Association* **106 (494)**, 746–762.

- [16] KRATZ, M., LOK, Y., AND MCNEIL, A. (2016). A multinomial test to discriminate between models. *ASTIN 2016 Proceedings, available online*.
- [17] KRATZ, M., LOK, Y., AND MCNEIL, A. (2016). Multinomial VaR Backtests: A simple implicit approach to backtesting expected shortfall. *ESSEC Working Paper 1617 & arXiv1611.04851v1*
- [18] TASCHE, D. (2002). Expected shortfall and beyond. *Journal of Banking and Finance* **26**, 1519–1533.
- [19] ZIEGEL, J. F. (2016). Coherence and elicibility. *Math. Finance* **26**, 901–918.

AVENUE BERNARD HIRSCH BP 50105
CERGY-PONTOISE, 95021 CEDEX, FRANCE
E-MAIL: kratz@essec.edu