

Document:	EC 2020/109/W.P.4
Agenda:	5
Date:	22 May 2020
Distribution:	Public
Original:	English

E



Investing in rural people

IFAD10 Impact Assessment Sensitivity Analyses and Implications for IFAD11

Note to Evaluation Committee

Focal points:

Technical questions:

Sara Savastano

Director
Research and Impact Assessment Division
Tel.: +39 06 5459 2155
e-mail: s.savastano@ifad.org

Alessandra Garbero

Senior Economist
Tel.: +39 06 5459 2458
e-mail: a.garbero@ifad.org

Dispatch of documentation:

Deirdre Mc Grenra

Chief
Institutional Governance and
Member Relations
Tel.: +39 06 5459 2374
e-mail: gb@ifad.org

Evaluation Committee — 109th Session
Rome, 19 June 2020

For: Information

I. Background

1. At its 127th session, the Executive Board reviewed the Consultation on the Tenth Replenishment of IFAD's Resources (IFAD10) Impact Assessment Report together with comments from the Independent Office of Evaluation of IFAD (IOE). The Board thanked Management for undertaking the assessment and welcomed the results. Further, the Board noted the need for Management to explore limitations in the current methodology and to work on improving the methodology going forward. Specifically, the Board asked Management to conduct a peer review of the methodology and further strengthen it with the support of an external expert, and also to consider sharing the sampling and methodology for review prior to undertaking the IFAD11 Impact Assessment (IA) exercise.
2. To fulfil these requests, Management hired an expert¹ for two purposes: to assess the methodology used for IFAD10 IA reporting to determine if there was any bias in the results generated by the selection of projects for the sample; and, based on this assessment, to confirm the selection approach for the IFAD11 IA.
3. Given the extent and complexity of the IFAD10 IA data, detailed sensitivity analyses were undertaken of the IA approach. The results of the analyses and the validation of the IFAD11 sample are presented in the appendix of this document.
4. IFAD conducted the IFAD10 selection using the protocol approved by the Board as part of the Development Effectiveness Framework.² As later suggested by the Board, the selection was complemented by a sensitivity analysis to test the robustness of the sample. The sensitivity analysis proved that any potential bias in the selection was negligible, indicating that the results presented in the IFAD10 Impact Assessment Report were valid. The added value of the corporate reporting methodology outweighs the bias encountered and should not undermine the overall undertaking that makes IFAD unique in terms of corporate reporting. The same methodology will be adopted to validate the IFAD11 sample. The validity of IFAD's approach to measure corporate reporting was confirmed.
5. IOE and the Executive Board were correct in raising the possibility of bias in the sample of selected projects. This is a legitimate concern in any sample, but particularly when it is impossible to undertake a random sample for the purpose of selection. A key lesson emerging from the sensitivity analyses is the need to consider potential bias at the project selection stage. Given that the IFAD11 selection has followed the protocol accepted by the Board, the results will be further validated through a sensitivity analysis. Alternative methodologies will be explored for IFAD12.
6. The findings of the analyses demonstrate that the methodology used for IFAD10 is valid. Furthermore, the sample selection process – which follows the protocol of the Development Effectiveness Framework – is reasonable for future IA activities. The approach does not imply any reputational risk for IFAD.

II. Summary of findings

7. A number of systematic analyses were conducted to detect the possible presence of bias in the results of the corporate IA. The possible bias related to the IFAD10 project sample selected and the extent to which those projects are representative of the portfolio of projects completing during IFAD10. IFAD has been conducting IAs with the goal of demonstrating both accountability and learning, hence representativeness, rigour and transparency are key principles of the methodology.

¹ Stefano Gagliarducci, Professor of Economics, University of Rome Tor Vergata, Department of Economics and Finance, and Research Affiliate, Einaudi Institute for Economics and Finance. Professor Gagliarducci has published in leading economics journals and has previously conducted research on publication bias.

² IFAD Development Effectiveness Framework (DEF) – Executive Board December 2016 (paragraph 58).

8. The findings of the analyses can be summarized as follows:

- (i) **Potential bias in the selection of projects for the IFAD10 IA is negligible.** A systematic assessment of the sample underlying the IFAD10 IA was conducted to ascertain the possible presence of selection bias, and the nature, direction and magnitude of such bias. To this end, all possible variables that might have influenced the selection were scrutinized. The variables included implementation performance ratings, i.e. the ratings available at the time of selection (July 2016). These are the only variables that could have influenced selection. Based on a total of 107 IFAD10 completed projects and across 24 implementation ratings, differences in average ratings between the sample of 15 per cent projects (equivalent to 19 projects) and the rest of the portfolio (88 projects) were tested for significance. It was found that bias was absent for the large majority of the variables and that such difference in average ratings was significant for just two variables: (i) disbursement rate and (ii) counterparts funds. An in-depth validation for the possible presence of selection bias across these two variables was therefore conducted.

To this end, subgroup meta-analyses were carried out to assess whether the magnitude of impact, as measured in the project-level IAs, was associated with the implementation performance rating class (from 1 to 6) to which a project belonged. This was instrumental in verifying whether there was an association between impact magnitude and rating scales. In other words, one would expect that projects with satisfactory ratings in terms of disbursement performance or counterpart funds would exhibit higher impact magnitudes compared with projects rated moderately satisfactory or unsatisfactory. This might have implied a positive relationship between performance and impact. However, Management's findings showed the absence of a clear relationship between the ratings class and the impact estimates, particularly in the case of disbursement performance ratings.

The market access strategic objective is a case in point: projects rated satisfactory or better for disbursement performance (namely coded 1 to 3) displayed the lowest impact (57 per cent) compared with moderately satisfactory ones (80 per cent) and projects rated unsatisfactory which, in turn, displayed the highest impact (89 per cent). This was also largely true for the other strategic objectives. Looking at the productive capacities strategic objective, and the other rating variable of counterpart funds, projects rated moderately satisfactory (namely coded 4) presented an 18 per cent impact compared with unsatisfactory ones (coded 1 to 3). Based on this assessment, Management concluded that the direction of the ratings – stronger performance – did not correspond to greater impact.

- (ii) **Selection bias corrections found that bias was marginal, indicating that the results presented in the IFAD10 Impact Assessment Report were valid.** As a second step, to further validate these results, two other methods were used to assess the need to adjust the results of the meta-analyses for the possible presence of selection bias. The first – selection bias correction using the Heckman model – was used to compute the project-specific probability of being selected for an IFAD10 IA and, in the case of selection, to adjust the corporate impact estimates for this bias. This probability takes into account selection drivers such as disbursement performance, counterpart funds, and other key observable variables.³The second "trim-and-fill"⁴ method draws on the meta-analysis literature and was initially considered for verifying the presence of bias generated by

³ Such an approach was combined with meta-regression and meta-analysis.

⁴ The "trim-and-fill method" is a popular tool to detect and adjust for publication bias.

deliberately including certain projects/studies in meta-analysis. A known limitation of trim-and-fill is that it can correct for publication bias that does not exist, underestimating effect sizes (the results, which are of low reliability, are provided in the appendix).⁵ Under the first method, it was found that selection bias was only marginally significant (at 10 per cent) for the market access strategic objective domain. Corrected impact estimates in this case using the Heckman method indicated that, *ceteris paribus*, impact was overestimated by less than 15 per cent.⁶

- (iii) **The projects selected for the IFAD11 IA sample will not lead to bias.** Last, a validation of the IFAD11 sample was conducted using the same approach. In light of the IFAD10 sensitivity analyses, a similar methodology was implemented to validate the IFAD11 sample of IAs and assess the presence of selection bias at the time of project selection (July 2018), using implementation performance ratings and other features of the portfolio universe. As part of the IFAD11 IA agenda, 24 out of 112 projects were selected for rigorous IA, representing 21.4 per cent of the portfolio, 20.9 per cent of total financing and 25.6 per cent of IFAD financing. As done for the validation of IFAD10, 24 implementation ratings as well as a number of other objective features of the portfolio (number of beneficiaries and financing variables for instance) were scrutinized for the possible presence of bias at the time of selection. No statistically significant differences either in objective features or in average implementation ratings were found across the sample of projects chosen (24) and the remainder of the projects in the universe (88 projects closing during IFAD11). Only the performance of monitoring and evaluation systems was statistically significant. These results demonstrate the absence of selection bias also in the case of the IFAD11 sample.
- (iv) In conclusion, in light of these statistical validations and sensitivity analyses, Management can confirm that there is no selection bias in either the IFAD10 or the IFAD11 samples of projects selected for corporate IA.

⁵ Terrin N, Schmid CH, Lau J, Olkin I (2003) Adjusting for publication bias in the presence of heterogeneity. *Statistics in Medicine* 22: 2113–2126.

⁶ 2 per cent in case of economic mobility, 15 per cent in case of market access, 10 per cent in case of production, and 6 per cent in the case of resilience. Nutrition results remain unchanged.

Appendix: Peer review of IFAD10 Impact Assessment Methodology

Stefano Gagliarducci, University of Rome Tor Vergata and EIEF

In collaboration with Alessandra Garbero and Sara Savastano, IFAD

Contents

1. Introduction	2
2. Background	2
3. Definitions.....	4
3.1. Representativeness	4
3.2. Sampling/selection bias	5
4. Descriptive analysis and selection bias.....	6
5. Strategies for addressing the selection bias	15
5.1. Modelling selection bias for impact assessment using observables features .	15
5.2. Publication bias.....	15
6. Results from the Sensitivity Analyses.....	16
6.1. Subgroup meta-analyses	16
6.2. Sample selection bias correction á la Heckman	16
7. Conclusions on IFAD10	19
8. Implications for IFAD11	20
8.1. IFAD11 Sample Validation.....	20
8.2. Conclusions for IFAD11	35
Annex I	
Annex II	
Annex III	
Annex IV	

1. Introduction

During the discussion of the IFAD10 Impact Assessment Report at the 127th Executive Board meeting on September 11, 2019, the Board recommended to conduct Sensitivity analyses to assess the robustness of the corporate impact estimates and verify the results. This recommendation was made in light of the comments received by the Independent Office of Evaluation (IOE) and other Stakeholders, indicating the possible presence of a bias in the meta-analysis estimates and projections, raising concerns around the credibility of the findings. Such bias concerned the choice of the IFAD10 sample of projects selected for an Impact Assessment (IA) and notably, that such projects are, according to IOE, not representative of the portfolio of projects completing during IFAD10.

IOE's argument was based on a descriptive analysis of the performance ratings⁷ at completion (or project completion reports ratings, in brief PCR). Their conclusion was that the projects selected for an impact assessment during IFAD10, seemed to include a large percentage of higher performing ones, therefore "potentially" yielding "biased" estimates of impact and possibly implying an overly optimistic vision of IFAD10 aggregate impact performance.

Systematic sensitivity analyses were therefore conducted to assess whether bias existed, and then investigate its magnitude. This is justified on transparency grounds, and on the fact that IFAD strongly believes in demonstrating accountability and learning, through rigorous methods.

Broadly speaking, sensitivity analysis is a process that allows the analyst to prove that the findings from a meta-analysis are not dependent on arbitrary or unclear decisions. In practice, they are aimed at repeating the meta-analysis, substituting alternative decisions or ranges of values for decisions that were arbitrary or unclear. For example, if the eligibility of some studies in the meta-analysis is dubious because they do not contain full details or are not representative, sensitivity analysis may involve undertaking the meta-analysis twice: first, including all studies and second, only including those that are definitely known to be eligible/representative. In this context and, through a weighting procedure, sensitivity analysis address the robustness of the results to the explicit inclusion of selection bias into the estimates, whereby this bias is assumed to be originated by the inclusion of a large number of projects with high performance ratings at completion, in the sample of the IFAD10 IAs.

The document presents the results of these analyses and is structured as follows. Section 3 recapitulates the background of IFAD approach to corporate reporting as stated in the Development Effectiveness Framework. Section 4 first introduce some definitions, section 5 presents a descriptive analysis, section 6 a literature review on the possible strategies to address bias, section 7 the results from the sensitivity analyses, section 8 concludes on IFAD10 and section 9 presents some implications for IFAD 11 and the corresponding validation of the IFAD11 sample of impact assessments.

2. Background

IFAD carries out project-level impact assessments (IAs) on a selection of projects (about 15 per cent) that are representative of the portfolio, to be able to measure corporate impact or aggregate development effectiveness. The latter requires a methodology that can attribute IFAD impact at the corporate level, e.g. provide an estimate of aggregate impact for the corporate indicators laid out in the IFAD Strategic Framework 2016-2025. The approach used is systematic, comprehensive, transparent, and builds upon the

⁷ Since 2005, in line with the practice adopted in many other International Financial Institutions (IFIs) and United Nations organizations, IOE uses a six-point subjective rating system (where 6 is the highest score and 1 the lowest score) to evaluate projects. In addition to reporting on performance based on the six-point rating scale, in 2007 IOE introduced the broad categories of "satisfactory" (rating coded 4 to 6) and "unsatisfactory" (rating coded 1 to 3) for reporting on performance across the various evaluation criteria.

IFAD9 Impact Assessment Initiative methodology as well as the IFAD10 Development Effectiveness Framework.

IFAD's Development Effectiveness Framework (DEF)⁸, approved by the Board in September 2016 lays out the selection protocol to assess projects suitability to undergo an impact assessment, specifying to the following criteria:

- (i) potential to learn lessons;
- (ii) feasibility of conducting a scientifically rigorous impact assessment;
- (iii) buy-in from the government and IFAD;
- (iv) the capacity of a project to represent IFAD's portfolio and
- (v) the relevance of the impact assessment for subsequent project phases.

A key factor of impact assessment, in addition to accountability, is learning; and learning needs to inform the design of new projects in the same country or elsewhere. This provides a public good for policymakers. Therefore, a major recommendation approved by the Board - in the Development Effectiveness Framework -- stated that "impact assessments should have been selected and structured to facilitate and maximize learning while recognizing the need for corporate reporting, and that an impact assessment agenda should be a multi-stakeholder and participatory process to ensure relevance" (IFAD, 2016⁹ pag.1).

Consequently, projects selected for IFAD10 IAs had to both display the potential for learning (innovative approaches or a clear evidence gap), while maintaining feasibility and have buy-in from the government.

In order to allow for adherence to the IFAD10 selection protocol, a working group was created to ensure that the selected projects were representative of the portfolio and revealed gaps for additional assessments, with a view to gaining an understanding of how projects fit into the portfolio. The expectation was that selected projects would have ultimately reflected the thematic and regional coverage of IFAD projects.

This led to a participatory process, finalized in September 2016, whereby projects selected for impact assessments were chosen in collaboration with IFAD's regional divisions to maximise this learning criteria. The divisions provided a list of projects suitable for inclusion based on the criteria specified according to the selection protocol. Subsequently, an appraisal was done to determine the impact assessments' feasibility in consultation with the regional divisions and relevant country directors.

Concerning corporate-level impact, IFAD's methodology to estimate aggregate development effectiveness involves a two-steps procedure whereby a meta-analysis of individual project-level impact assessment estimates is conducted in the first stage to compute aggregate corporate impacts, and a projection is conducted in the second stage to extrapolate impacts to the rest of the portfolio and estimate number of people benefiting across the portfolio¹⁰.

⁸ The DEF was developed based on the lessons learned from the experience in demonstrating impact as part of the IFAD9 Impact Assessment Initiative. See [EB 2016/119/R.12](#)

⁹ International Fund for Agricultural Development (IFAD), 2016. Development Effectiveness Framework ([EB 2016/119/R.12](#)).

¹⁰ As far as the projection approach is concerned, this refers to a methodology that allows the estimated impact to be extrapolated to the whole IFAD portfolio, in order to obtain an assessment of the number of people that have benefited from IFAD investments. The corporate impact is interpreted as percentage change gain in each of the Strategic Objectives (SOs) and on IFAD's overarching goal. To translate this into the number of beneficiaries who benefited from IFAD's investments, distributional assumptions are needed to extrapolate the corporate estimates to the universe of beneficiaries in the portfolio.. The IFAD10 projection universe includes 107 projects, and is defined as the total number of projects completing during the replenishment period (2016-2018). As the projection require estimates of beneficiaries reached across the whole universe, the additional challenge has been to aggregate the number of beneficiaries for the overall portfolio. The information on the number of beneficiaries in the IFAD10 portfolio can be extracted from project documentation and IFAD internal reporting systems. Projected beneficiaries impacted are calculated based on the number of actual beneficiaries belonging to the universe of 107 projects. The latter amount to around 65.3 million beneficiaries. At the basis of the extrapolation, there are two main assumptions. One concerns the distribution of impacts, where the assumption is that corporate impacts are normally distributed with means and standard errors corresponding to the ones estimated empirically while obtaining aggregate impact estimates from the 17 impact studies covering 19 projects (equivalent to 18 per cent of the universe, actually). The second

The aggregation is systematically done via a meta-analysis, a statistical procedure for combining data from multiple studies. Meta-analysis was pioneered in medical studies in the late seventies and then exponentially applied to clinical research. The meta-analysis is a study design used to systematically assess previous research studies to derive conclusions about a specific drug/treatment/research (or in our case policy) question. Outcomes from a meta-analysis may include a more precise estimate of the effect of treatment or risk factor for disease, or other outcomes (Haidich, 2010)¹¹. More broadly, meta-analysis is defined as “the statistical analysis of a large collection of results for the purpose of integrating the findings” (Glass 1976¹²). In other words, it is “a quantitative summary of statistical indicators reported in similar empirical studies” (Brander et al. 2006¹³).

In the context of IFAD10 IA, the meta-analysis is a statistical procedure that aggregate the results of the 15% of projects on which an individual study is conducted. The outcome of the analysis is a proxy for an average effect (the treatment effect, or effect size) of the impact of IFAD’s overall portfolio. Once aggregated, corporate impacts were computed as percentage changes over the comparison group for each Strategic Objective (SO), notably production, market access or participation, and resilience, and for the overall IFAD goal of increased economic mobility.

3. Definitions

In this section, some definitions that are going to be useful for an understanding of the remainder of the document are provided.

3.1 Representativeness

First, the concepts of representativeness, population or universe, and sample is defined: a representative sample is one that matches some characteristic of the underlying population, usually the characteristic that one is targeting with the research. In the context of IFAD10, the population refers to the population of projects that are in the universe of projects completing during IFAD10 (around 107 projects). The sample under analysis is defined as the 19 projects chosen for an impact assessment study.

As mentioned before, the IFAD10 selection protocol for the 19 IFAD10 Projects was based on a number of criteria to ensure representativeness of the portfolio. To what extent such criteria were as good as random is open to question. They could be quasi-random, in the sense that ex-ante, or at the time of selection, it was not possible to ascertain all of them.

In practice, it is almost impossible to ensure randomness due to a number of factors, such as feasibility of the impact assessment itself, knowledge asymmetries, political considerations and stakeholders buy-in, among others. It is worth recalling that the issue

assumption is about defining what benefiting means in terms of exceeding a certain threshold. The projected number of beneficiaries impacted by IFAD’s investments can be obtained by setting a threshold of at least 20 per cent for impact gains. Using estimates on the aggregate impacts and knowledge of the portfolio, one can then obtain projected number of beneficiaries benefiting above a 20 per cent threshold. In summary, projected beneficiaries impacted are obtained by randomly drawing a normal distribution of impacts with means and standard errors centred to the ones empirically estimated from aggregate impact distributions, thereby assuming that benefits are randomly and normally distributed and are above a specific threshold.

¹¹ “Important medical questions are typically studied more than once, often by different research teams in different locations. In many instances, the results of these multiple small studies of an issue are diverse and conflicting, which makes the clinical decision-making difficult. The need to arrive at decisions affecting clinical practice fostered the momentum toward “evidence-based medicine. Evidence-based medicine may be defined as the systematic, quantitative, preferentially experimental approach to obtaining and using medical information.” Haidich AB, 2010, *Meta-analysis in medical research*. Hippokratia. 2010 Dec.14 (Suppl 1):29-37.

¹² Glass, G. (1976). Primary, Secondary, and Meta-Analysis of Research. *Educational Researcher*, 5(10), 3-8.

¹³ Brander L.M., et al. (2007). The recreational value of coral reefs: A meta-analysis, *Ecological Economics*, Vol. 63, Issue 1, 2007, 209-218.

of representativeness of the impact assessment sample was also raised during the previous replenishment cycle (IFAD9). Maintaining the integrity of the random selection conducted during IFAD9, was extremely difficult due to the above-mentioned factors. In that instance, projects were selected according to a number of criteria: 1) feasibility (suitable for an ex post impact assessment); 2) with the overarching aim of measuring the poverty reduction impact and, 3) statistically representative of the portfolio of activities undertaken by IFAD during IFAD9.

Therefore a representative sample of projects to be evaluated was determined by drawing a stratified random sample (a total of 41 projects, i.e. 26 first-choices and 15 reserves) from the universe of projects (with available datasets) closing between 2010 and 2015.

However, maintaining the integrity of the random sample proved difficult, as some randomly selected ones had to be replaced owing to both political and practical concerns (conflict setting, absence of PMU or key informants essential to gather retrospective information about projects, and impossibility to determine a counterfactual, among others). An internal consultation (in 2012) with IFAD Regional Directors and divisional representatives was then conducted to endorse the list of randomly selected projects to be evaluated by the external research partners. This process led to the replacement of 11 randomly selected projects with a set of purposively selected ones (purposive evaluations), given their strategic relevance and overall performance across the portfolio. Two of the purposively selected projects were dropped (namely, those in India and Senegal) after discussing the feasibility with internal staff. This last factor showed that even for “cherry picked” projects feasibility of an impact assessment was not guaranteed.

Notwithstanding these issues, IFAD sought to maintain the integrity of the representative sample and decided to maintain the randomly selected projects excluded from the final list of ex post evaluations and conduct the ex post assessments in-house with secondary datasets (14 Shallow dives).

Regarding IFAD10, and as noted above, a selection protocol was followed to ensure representativeness of the portfolio. The rationale for using a protocol is similar to what is normally conducted in the medical field, which is to randomly assign patients into treatment and control groups. As such, these protocols have features of quasi-randomness – as patients are selected into treatment – across a population of eligible patients¹⁴.

3.2 Sampling/selection bias

Selection bias is problematic because it is possible that a statistic computed of the sample is systematically erroneous. Selection bias can lead to a systematic over- or under-estimation of the corresponding parameter in the population. Selection bias occurs in practice as it is practically impossible to ensure perfect randomness in sampling (see before). If the degree of misrepresentation is small, then the sample can be treated as a reasonable approximation to a random sample. Also, if the sample does not differ markedly in the quantity being measured, then a biased sample can still be a reasonable estimate.

Selection bias is mostly classified as a subtype of selection bias, sometimes specifically termed sample selection bias, but some classify it as a separate type of bias. A distinction,

¹⁴ There is currently a heated debate around the topic of randomized controlled trials and whether they should be generally considered the gold standard, namely the best method to infer causality. It is worth recalling that in medical studies, researchers often choose not to randomize the intervention for one or more of the following reasons: (1) ethical considerations, (2) difficulty of randomizing subjects, (3) difficulty to randomize by locations (by region in the case of IFAD portfolio), (4) small available sample size (Harris et al. 2006).

albeit not universally accepted, of selection bias is that it undermines the external validity of a test (the ability of its results to be generalized to the entire population), while selection bias mainly addresses internal validity for differences or similarities found in the sample at hand. In this sense, errors occurring in the process of gathering the sample or cohort cause selection bias, while errors in any process thereafter cause selection bias. In this specific case, it refers to selection bias.

However, selection bias and selection bias are often used synonymously.

4 Descriptive analysis and selection bias

In this section, a first assessment of the presence of selection bias is provided by presenting descriptive statistics that characterize the universe of IFAD10 projects, compared to the sample chosen for impact assessments. These descriptives are essential to understand the extent and the severity of the bias based on observable features. In presence of a selection bias, one should expect the two groups to differ significantly over an array of observable dimensions.

Recall that the sample is made of 19 projects, and that the universe is composed of 107 IFAD10 projects slated to close during IFAD10 at the time of selection. After projects were selected for assessment, a subset of the selected projects were extended such that their closing dates are now in IFAD11¹⁵.

The main argument against lack of representativeness cited by IOE was that the sample of IFAD10 IAs included a large majority of high performing projects as displayed by IOE's analyses of performance indicators at completion (PCR ratings).

Thus, Table 1 reproduces the one presented by IOE in their summary document for the Evaluation Committee Session (EC) held in September 2019. Average performance indicators at completion (PCR ratings) are displayed for the sample of 19 projects evaluated as part of IFAD10 IAs and the remaining 88 projects not evaluated out of the total universe of 107 projects.

PCR ratings at completion are subjective ratings with a six-point measurement scale system ranging from 6 to 1, with 1 being the lowest score across each criterion. At the time of IFAD10 selection, when the sample of 15% of projects was identified, none of these scores were available for consultation nor officially available within IFAD official system¹⁶.

While comparing the two tables, a number of issues became apparent. The first, is the lack of definition of the universe of projects analysed e.g. the total number of observations (projects), in IOE's document. As a consequence statistics for unselected projects completing during IFAD10 (columns 3 and 4) varied across some indicators, while results for the IFAD10 sample (namely columns 1 and 2) coincided with IOE calculations. Also, PCR scores were not available for all the unselected projects, therefore the total number

¹⁵ Notably Bangladesh (CCRIP), Kenya (SCDP), Sao Tome (PAPAC), Rwanda (PRICE) will now close in IFAD11.

¹⁶ Specifically, regarding features that might have driven the IFAD10 IAs selection process, and alter the representativeness of the sample of the IFAD10 projects portfolio, the following ones were available at the time of the selection, notably in 2016: the project type or sector, the region of implementation, the size of outreach, the disbursement performance, and the implementation performance indicators. As projects were ongoing, project completion report ratings (the one verified by IOE) were not available to inform the selection.

of observations by indicator varies between 64 to 88 projects (column 3). Last, and similar to IOE's table, final PCR ratings are only available for 13 of the 19 projects that underwent an IA. The unavailability of the PCR ratings for the whole IA sample, is due to the fact that PCRs cannot be finalized if projects completion dates are extended. This was the case of 6 projects out of 19, whereby their completion dates were extended into IFAD11.

Therefore, Table 1 shows the difference across PCR ratings as presented by IOE in their comments. T-tests were run for the statistical significance of the difference in means (balance tests).

Before commenting the table, it is important to highlight that, while this is certainly an informative exercise, the latter should be taken with caution. As stated in the DEF, ex-post impact assessment should ideally occur prior to the closure of the project, so project completion reports can benefit from the impact assessment findings. If so, PCRs ratings incorporate IA findings when available – hence potentially influencing the direction of the final rating. Therefore, from a statistical standpoint, PCR ratings should not be used to assess the presence of selection bias, as they are positively affected by the mere virtue of a project being under evaluation.

Table 1: Balance tests: PCR ratings

	Average PCR ratings (IFAD10 IA sample)		Average PCR ratings (completing IFAD10 projects 2016-2018)		Sample - Unselected	
	(1) N. projects	(2) Mean	(3) N. projects	(4) Mean	(5) Diff. in Means	(6) P-score
Relevance	13	5.2	75	4.6	0.6	0.005
Effectiveness	13	4.8	76	4.2	0.6	0.005
Efficiency	13	4.4	76	3.8	0.6	0.026
Sustainability	13	4.4	76	3.8	0.5	0.017
Project performance	13	4.7	75	4.1	0.6	0.002
Rural poverty impact	13	4.8	75	4.1	0.7	0.001
Gender equality and women's empowerment	13	4.6	88	0.6	0.2	0.38
Innovation	13	4.8	76	4.4	0.4	0.111
Scaling up	13	4.8	75	4.4	0.5	0.066
Environment and natural resource management	13	4.5	73	4.1	0.4	0.056
Adaptation to Climate Change	11	4.5	64	4.1	0.5	0.022
IFAD performance	13	4.8	76	4.3	0.5	0.01
Government performance	13	4.7	76	4.1	0.6	0.014
Overall project achievement	13	4.8	75	4.2	0.6	0.004

Source: Calculations based on IFAD10 IA sample and data extracted from IOE ratings database.

Table 1 shows average subjective rating scores across 14 mandatory criteria¹⁷, used by IOE to evaluate projects at completion. However, means of selected and unselected projects are based on the universe of projects as defined by Management in the IFAD10 Report (107 completed projects).

Given the above concern with the use of PCR ratings, in what follows, the significance of differences in pre-determined characteristics is tested. These essentially are baseline characteristics and include objective features and 24 implementation ratings as measured at the beginning of the project (i.e., the first indicator of performance that is available in the system). Implementation ratings are monitored during the lifespan of the project. These are the ones that, effectively, should have informed projects' selection at the beginning of the IA process. Note that while the first three indicators in the table are objective, notably project duration, number of beneficiaries and total approved funding, performance indicators are self-assessed and are expressed on a rating scale (1-6) ranging from unsatisfactory, to highly satisfactory¹⁸.

Zooming in, note how the projects selected for impact assessments were similar on average in terms of financing and number of actual beneficiaries to the universe of projects. The average approved financing across the sample of IAs was \$51.7 million, and the average in the universe was of \$50.9 million. In terms of beneficiaries, the average number of beneficiaries was 610,556 in the universe and 490,339 in the IA sample, but this difference is not statistically significant. In almost all performance ratings categories, the IAs performed slightly better than the universe of projects, on average. However it is

¹⁷ Based on IOE Manual (2015) pp. 38 -40. These definitions build on the OECD/DAC Glossary of Key Terms in Evaluation and Results-Based Management; the Methodological Framework for Project Evaluation agreed with the Evaluation Committee in September 2003; the first edition of the Evaluation Manual discussed with the Evaluation Committee in December 2008; and further discussions with the Evaluation Committee in November 2010 on IOE's evaluation criteria and key questions. Rural poverty impact is defined as the changes that have occurred or are expected to occur in the lives of the rural poor (whether positive or negative, direct or indirect, intended or unintended) as a result of development interventions. Project performance is an average of the ratings for relevance, effectiveness, efficiency and sustainability of benefits. Relevance measures the extent to which the objectives of a development intervention are consistent with beneficiaries' requirements, country needs, institutional priorities and partner and donor policies. It also entails an assessment of project design and coherence in achieving its objectives. An assessment should also be made of whether objectives and design address inequality, for example, by assessing the relevance of targeting strategies adopted. Effectiveness is the extent to which the development intervention's objectives were achieved, or are expected to be achieved, taking into account their relative importance. Efficiency is a measure of how economically resources/inputs (funds, expertise, time, etc.) are converted into results. Sustainability of benefits (or simply sustainability) is the likely continuation of net benefits from a development intervention beyond the phase of external funding support. It also includes an assessment of the likelihood that actual and anticipated results will be resilient to risks beyond the project's life. Gender equality and women's empowerment measures the extent to which IFAD interventions have contributed to better gender equality and women's empowerment, for example, in terms of women's access to and ownership of assets, resources and services; participation in decision making; work load balance and impact on women's incomes Nutrition and livelihoods. Innovation and scaling up (OR scaling up) measures the extent to which IFAD development interventions: (i) have introduced innovative approaches to rural poverty reduction; and (ii) have been (or are likely to be) scaled up by government authorities, donor organizations, the private sector and others agencies. Environment and natural resource management represents the extent to which IFAD development interventions contribute to resilient livelihoods and ecosystems. The focus is on the use and management of the natural environment, including natural resources defined as raw materials used for socio-economic and cultural purposes, and ecosystems and biodiversity – with the goods and services they provide. Adaptation to climate change is the contribution of the project to reducing the negative impacts of climate change through dedicated adaptation or risk reduction measures. Performance of Partners (IFAD and Government): This criterion assesses the contribution of partners to project design, execution, monitoring and reporting, supervision and implementation support, and evaluation. The performance of each partner will be assessed on an individual basis with a view to the partner's expected role and responsibility in the project life cycle. Finally, overall project achievement provides an overarching assessment of the intervention, drawing upon the analysis and ratings for rural poverty impact, relevance, effectiveness, efficiency, sustainability of benefits, gender equality and women's empowerment, innovation and scaling up, as well as environment and natural resources management, and adaptation to climate change.

¹⁸ Ratings of project performance should be consistent with the findings of progress reports and of the supervision mission report. By rating each indicator, different criteria are applied as explained below, however in general the ratings are:

(6) Highly satisfactory. Targets/requirements met or exceeded. Considered as best practice.

(5) Satisfactory. Targets/requirements met with only minor delays or set-backs.

(4) Moderately satisfactory. Most targets/ requirements met but delays or set-backs experienced.

(3) Moderately unsatisfactory. Some targets/ requirements met but issues/constraints have negatively affected implementation.

(2) Unsatisfactory. Few targets/requirements met. Issues/constraints remain unresolved. Delays have seriously undermined implementation.

(1) Highly unsatisfactory. Almost no targets/ requirements met. Consideration should be given to cancellation/suspension.

important to note that these differences are not statistically significant for the majority of indicators presented – except for the following ratings:

- Assessment of the Overall Implementation Performance* : significant at 10% level.
- Acceptable Disbursement Rate** : significant at 5% level.
- Counterparts Funds** : significant at 5% level.
- Coherence between AWPB and Implementation* : significant at 10% level.

Table 2: Balance tests: implementation performance ratings (baseline characteristics)

	<i>IFAD10 IA Sample</i>		<i>IFAD10 projects</i>	<i>Unselected</i>	<i>Sample - Unselected</i>	
	N. projects (1)	Mean (2)	N. projects (3)	Mean (4)	Diff. in Means (5)	P-score (6)
Project Duration	19	8.16	88	8.27	-0.12	0.863
Beneficiaries	19	490 339	88	636 512	-146 173	0.732
Approved Funding	19	51 712 292	88	50 700 000	986 098	0.957
Assessment of the Overall Implementation Performance	19	4.11	88	3.85	0.25	0.059
Likelihood of Achieving the Development Objective Effectiveness	19	4	88	3.98	0.02	0.857
	12	3.83	69	3.88	-0.05	0.74
Targeting and Outreach	19	4.26	88	4.11	0.15	0.224
Gender equality & women's participation	19	4.05	88	3.99	0.06	0.666
Agricultural Productivity	15	4.13	71	3.94	0.19	0.147
Adaptation to Climate Change	2	4	11	3.91	0.09	0.863
Institutions and Policy Engagement	18	4.11	78	4.01	0.1	0.549
Human and Social Capital and Empowerment	15	4.13	77	3.92	0.21	0.177
Quality of Beneficiary Participation	19	3.95	88	4.06	-0.11	0.401
Responsiveness of Service Providers	19	3.89	88	3.97	-0.07	0.592
Environment and Natural Resource Management	2	4	13	3.77	0.23	0.607
Exit Strategy	11	4.09	58	3.97	0.13	0.387

Potential for Scaling-up	14	4.21	72	4.07	0.15	0.366
Quality of Project Management Knowledge Management	19	3.95	88	3.85	0.1	0.631
	16	4.19	76	4.03	0.16	0.243
Coherence between AWPB and Implementation	17	4.12	81	3.79	0.33	0.064
Performance of M&E System	19	3.74	87	3.83	-0.09	0.574
Acceptable Disbursement Rate	19	4.21	88	3.43	0.78	0.028
Quality of Financial Management	17	4.12	79	3.9	0.22	0.218
Quality and Timeliness of Audit	19	4.11	87	4.01	0.09	0.582
Counterparts Funds	19	4.42	88	4.01	0.41	0.031
Compliance with Loan Covenants	19	4.21	88	4.02	0.19	0.189
Procurement	19	4.16	88	4	0.16	0.292

Table 3 and Table 4 show the distribution of the sample of IAs projects by IFAD's region and project sector or type. In the universe, 30 projects were in the Asia and Pacific Region (APR) followed by 26 in Western and Central Africa (WCA), 20 in Eastern and Southern Africa (ESA), 18 in Latin America and Caribbean (LAC), and 13 in the North East and Northern Africa region (NEN). As far as the IAs Projects' Sample is concerned, the majority of IAs (six) were conducted in ESA, while five were conducted in APR, four in WCA, three in LAC, and one in NEN. Table 22 in the Annex presents the mean performance by region and shows similar results i.e. that none of the mean ratings are statistically different across the IA sample and the unselected projects, although there is more variation, largely due to the lower number of overall projects with each region.

Turning to the project sector or type, a variable that is quite broad in the current classification system, the majority of projects in the universe are classified as agricultural development (37), rural development (34), and credit (14). However, no credit projects were selected for assessment in IFAD10 and over 40% of all IAs were of rural development projects. Nevertheless, because the project sector categorization is extremely broad, contains considerable overlap among categories, and is insufficiently informative about the true nature of the project, it lacks utility for the conduction of rigorous sensitivity analysis or bias estimation.

Table 3: Distribution of Projects in the Universe and in the IAs sample by Region

Universe by Region			IAs by Region		
BU	Projects	%	BU	Projects	%
APR	30	28.04	APR	5	26.32
ESA	20	18.69	ESA	6	31.58
LAC	18	16.82	LAC	3	15.79
NEN	13	12.15	NEN	1	5.26
WCA	26	24.30	WCA	4	21.05
Total	107	100.00	Total	19	100.00

Table 4: Distribution of Projects in the Universe and in the IAs sample by Sector or Project type

Universe by Project			IAs by Project		
Sector	Projects	%	Sector	Projects	%
AGRIC	37	34.58	AGRIC	6	31.58
CREDI	14	13.08	CREDI	0	0.00
FISH	2	1.87	FISH	0	0.00
IRRIG	7	6.54	IRRIG	1	5.26
LIVST	4	3.74	LIVST	2	10.53
MRKTG	6	5.61	MRKTG	1	5.26
RSRCH	3	2.80	RSRCH	1	5.26
RURAL	34	31.78	RURAL	8	42.11
Total	107	100.00	Total	19	100.00

Finally, implementation performance ratings were combined (Table 5) to assess and test for differences across proportions/percentage of projects rated satisfactory both in the IA sample(19) and in the universe of projects completing during IFAD10 (107). Note that there was not much variation in project scores in either the universe or the IA sample with

the highest density of projects around scores of four and five out of six. Specifically, when the indicators are transformed to indicate whether a project scored a satisfactory (4-6) or unsatisfactory (1-3) rating, it is apparent that in both the IFAD10 IA sample and the universe the majority of projects received satisfactory ratings. In the IA sample, between 84 and 100 percent of projects received satisfactory ratings and in the universe between 75 and 98 percent of projects did.

Although there are differences in the relative frequency of unsuccessful projects, it is clear that the majority of portfolio projects receive satisfactory scores. As such, it is reasonable that a high proportion of the sample would be high performing projects.

In summary, a conclusion that can be drawn from this analysis, is that bias is absent for the majority of baseline indicators presented – except for the following two ratings:

- Acceptable Disbursement Rate**: significant at 5% level.
- Counterparts Funds** significant at 5% level.

Further analyses are therefore conducted on these variables in question in the following sections.

Table 5: Balance tests: Proportions of Projects Rated Satisfactory

	Average performance ratings	Average performance ratings	Difference in IAs and Universe Means	Proportion of IAs Rated	Proportion of Universe Rated	Sample-Universe	Significance
	IFAD10 IA Sample	(Completing IFAD10 projects 2016-2018 : universe(107))		Satisfactory	Satisfactory		
	Mean (1)	Mean (2)	Diff. in Means (3)	Proportion (4)	Proportion (5)	Diff. in Proportion (5)	P-score (6)
Assessment of the Overall Implementation Performance Likelihood of the Achieving Development Objective Effectiveness Targeting and Outreach Gender equality & women's participation Agricultural Productivity Adaptation to Climate Change Institutions and Policy Engagement Human and Social Capital and Empowerment Quality of Beneficiary Participation	4.11	3.9	0.21	100%	92%	8%	0.436
	4	3.98	0.02	100%	84%	16%	0.561
	3.83	3.88	-0.04	100%	83%	17%	0.957
	4.26	4.14	0.12	100%	98%	2%	0.336
	4.05	4	0.05	89%	89%	1%	0.93
	4.13	3.98	0.16	100%	84%	16%	0.569
	4	3.92	0.08	94%	90%	5%	0.685
	4.11	4.03	0.08	100%	82%	18%	0.537
	4.13	3.96	0.18	95%	87%	8%	0.483
	3.95	4.04	-0.09	100%	95%	5%	0.833

Responsiveness of Service Providers	3.89	3.95	-0.06	95%	87%	8%	0.221
Environment and Natural Resource Management	4	3.8	0.2	100%	87%	13%	0.582
Exit Strategy	4.09	3.99	0.11	84%	75%	9%	0.356
Potential for Scaling-up	4.21	4.09	0.12	100%	90%	10%	0.355
Quality of Project Management	3.95	3.87	0.08	84%	82%	2%	0.851
Knowledge Management	4.19	4.05	0.13	100%	85%	15%	0.293
Coherence between AWPB and Implementation	4.12	3.85	0.27	84%	69%	15%	0.317
Performance of M&E System	3.74	3.81	-0.07	78%	72%	7%	0.993
Acceptable Disbursement Rate	4.21	3.57	0.64	95%	65%	29%	0.025
Quality of Financial Management	4.12	3.94	0.18	89%	79%	11%	0.764
Quality and Timeliness of Audit	4.11	4.03	0.08	89%	93%	-3%	0.486
Counterparts Funds	4.42	4.08	0.34	79%	74%	5%	0.194
Compliance with Loan Covenants	4.21	4.06	0.15	100%	87%	13%	0.492
Procurement	4.16	4.03	0.13	89%	74%	16%	0.681

5 Strategies for addressing the selection bias

Given that there are statistically significant differences in only two observable features of the projects subject to an impact assessment, two main strategies are considered to assess the need for adjusting for the possible selection bias.

5.1 Modelling selection bias for impact assessment using observables features

Since some information on the universe is available (including implementation performance), the Heckman approach can be adapted (1979)¹⁹ to compute the likelihood of a project being selected for an assessment compared to the rest of the universe, conditional on the available observed characteristics. The meta-analysis is then run once again after reweighting each project by its probability of being selected into an IA.

The success of this approach, of course, rests on the number and on the importance of the observed factors in driving the selection into an IA. If one can plausibly argue that selection of IA projects depends mostly on the observable (rather than unobservable) characteristics that are observed, such as project type, region, financing and implementation performance ratings, the meta-analysis results can be adjusted for selection bias based on observables.

5.1 Publication bias

The second approach considered draws on the meta-analysis literature and treat our sampling issue as a classical publication bias problem. The latter refers to the distortion of meta-analysis outcomes due to the higher likelihood of publication of statistically significant studies rather than non-significant studies. This is similar to the problem at hand – where impact assessment estimates are only available for the projects evaluated, hence less performing projects are not observed – hypothetically – in the sample. Therefore, presenting this kind of sensitivity analysis would be useful here too.

In order to test for the presence or absence of publication bias, first, a funnel plot can be used. In essence, studies are plotted on a scatter plot with effect size on the x-axis and precision or total sample size on the y-axis. If the points form an upside-down funnel shape, with a broad base that narrows towards the top of the plot, this indicates the absence of a publication bias. On the other hand, if the plot shows an asymmetric shape, with no points on one side of the graph, then publication bias can be suspected. Second, to test publication bias statistically, Begg and Mazumdar's rank correlation test or Egger's test can be used. If publication bias is detected, the trim-and-fill method can be used to correct the bias (Shi et al, 2019).

A known limitation of Trim-and-Fill is that it can correct for publication bias that does not exist, underestimating effect sizes (Terrin et al, 2003). Results of this method are therefore optional and presented in Annex IV.

¹⁹ Heckman, J. (1979). Sample Selection Bias as a Specification Error. *Econometrica*, 47(1), 153-161.

6 Results from the Sensitivity Analyses

6.1 Subgroup meta-analyses

Subgroup meta-analyses are presented to assess whether the magnitude of the impact is associated with the rating class to which a project belongs, and verify whether there is a gradient between impact magnitude and rating scale. The rationale for this is due to the fact that in Section 4 it was shown how two baseline ratings were in fact unbalanced between the IFAD10 IA sample of projects and the rest of the universe.

Specifically, two features (Acceptable Disbursement Rate and Counterparts Funds) were significantly different at 5% level between projects evaluated and the unselected ones in the universe. Therefore, the extent of bias is investigated, notably whether there is a relationship between positive ratings and impacts. If bias existed, one would expect to see patterns or gradients such as the following, namely the higher the rating the higher the impact in the IFAD10 estimates of the aggregate effect sizes, as this is the main argument by IOE, notably that impact is overestimated due to higher performing projects.

Results across the Strategic Objectives (SOs), notably Market access, Resilience, Production as well as the cross-cutting theme Nutrition and the overarching IFAD goal of Economic mobility, are presented in the Annex (). Forest plots, the visual representations commonly used for meta-analytic results are employed for the purpose. The reader needs to focus on the diamonds, which represent the size of the effect – grouped according to ratings: unsatisfactory (below 3), moderately satisfactory (4), and satisfactory plus (ratings higher than 5). The “overall” diamond represents our pooled effect size, the one presented to the Board in September.

It is remarkable to see how there is no clear relationship between the ratings class and the impact estimates, particularly in the case of disbursement performance ratings. This is reassuring and corroborates the absence of bias in the impact estimates due to selection of higher performing projects.

For example, note how in the case of the market access SO domain, projects rated unsatisfactory in terms of disbursement performance ratings, display the highest impact (1.89 equivalent to 89%), compared with satisfactory plus (57%) and moderately satisfactory ones (80%). This is also largely true for the other SOs.

Also in the case of the performance rating “counterpart funds”, note how projects rated as unsatisfactory have a higher impact (29 %) compared with the moderately satisfactory ones (18%), in the domain of Production. Tables are presented in Annex II.

Therefore, a major conclusion that can be drawn from this analysis is that the direction of the ratings e.g. stronger performance, does not correspond to a larger impact.

6.2 Sample selection bias correction á la Heckman

In this section, one of the sensitivity analyses described in Section 6, palatable for adjusting for the bias, is presented. The second method – trim and fill – adjusts for bias non-parametrically and given its low reliability, is presented in the Annex.

In the approach presented here, a three-stage estimation procedure is applied, whereby the results are corrected for the presence of observable and unobservable selection using an approach a la Heckman (1979), combined with meta-regression and meta-analysis.

In the first stage, a probit regression model is estimated over the universe of projects on the variables determining the selection, notably the ones that are assumed as observable drivers for the selection into IA (project sector, region, and the significant performance ratings). These estimates are used to obtain an Inverse Mills Ratio (IMR), and introduce the latter in a meta-regression in order to estimate corrected standard errors. The variance associated with this corrected standard error is then computed and summed to the original variance, to obtain final “corrected standard errors” (second stage). These final corrected standard errors are then meta-analyzed along with the original effect sizes to derive a pooled effect size adjusted for bias (third stage).

It is important to note that IMR is only marginally significant (at 10% level) in the meta-regression pertaining to the market access SO domain. For the rest of SOs it is not significant - this implies that the hypothesis of selection bias is rejected in the case of these SOs, and that it weakly holds in the context of market access.

Results are summarized in Table 6, where “observed” refers to the original impact estimates presented in the official IFAD10 Report and “adjusted” refers to the ones corrected for sample selection bias. The full set of tables is presented in the Annex III.

Note how the results maintain the integrity of the baseline random-effect meta-analysis model - the one presented in the original IFAD10 Report - and show minimal discrepancies. Specifically, results based on this scenario also remain largely positive, and indicate that, *ceteris paribus*, impact is overestimated by 2% in case of economic mobility, 15% in case of market access, 10% in case of production, and 6% in the case of resilience. Nutrition results remain unchanged.

However, given the lack of significance of the IMR in the second stage regression, Management concludes that there is no selection bias in the corporate impact estimates and that bias adjusted estimates are not needed and represent an over-correction.

Table 6: Results from the Sample selection bias correction

Production				
	N. projects	ES	Lower CI	Upper CI
Observed	17	1.44	1.26	1.64
Adjusted	17	1.33	1.24	1.43
Market Access				
	N. projects	ES	Lower CI	Upper CI
Observed	16	1.76	1.45	2.14
Adjusted	16	1.51	1.32	1.72
Resilience				
	N. projects	ES	Lower CI	Upper CI
Observed	17	1.13	1.02	1.25
Adjusted	17	1.06	1.02	1.1
Nutrition				
	N. projects	ES	Lower CI	Upper CI
Observed	16	1.01	0.99	1.03

Adjusted	16	1.01	0.99	1.02
Economic Mobility				
	N. projects	ES	Lower CI	Upper CI
Observed	17	1.74	1.51	1.97
Adjusted	17	1.72	1.42	2.02

7 Conclusions on IFAD10

In this document, a number of sensitivity analyses are presented, to assess the presence, direction and magnitude of the possible selection bias inherent in the sampling of projects chosen to be evaluated under IFAD10.

Results highlight that the bias is absent and, if anything, negligible. Through a detailed descriptive analyses it is shown that almost all the pre-determined e.g. baseline features of the IFAD10 IA sample and of the ones of the unselected projects are largely balanced – e.g. they are not statistically different – with the only exception of a couple of implementation ratings. Upon further investigation, it was found that the direction of the ratings does not imply a larger estimated impact, allowing one to conclude that projects rated highly unsatisfactory on certain attributes exhibit higher effect sizes compared with satisfactory projects. This finding strongly hints that corrective actions are put in place by implementers across the project lifetime to influence ratings towards more positive ones, particularly at completion. This factor corroborates Management's choice of not employing ratings at completion (PCR ratings) for an assessment of selection bias as the latter are endogenous (e.g., influenced by the evaluation process) and may be inflated by many reasons, the first being that ratings may reflect corrective actions by implementers, and second, that ratings do incorporate the findings of the impact assessments when available.

This analysis is complemented by an assessment of the need to correct for sample selection bias. To this end, two approaches are considered, notably the sample selection bias correction a la Heckman and the trim-and-fill approach. These are meant to more formally assess the presence and the magnitude, respectively, of any possible sample selection bias.

Given that information about the observable factors that might influence selection are available in the system, a sample selection bias correction a la Heckman is the preferred approach and is applied in the meta-analytic context.

Results based on this scenario show that selection bias does not hold and it is weakly present only in the estimations of corporate impact for market access. After computing bias adjusted estimates – the latter remain largely positive, and indicate that, *ceteris paribus*, impact was overestimated by 2% in case of economic mobility, 15% in case of market access, 10% in case of production, and 6% in the case of resilience. Nutrition results remain unchanged.

The trim-and-fill method is instead a popular tool to detect and adjust for publication bias, in other words the bias originated by ad-hoc inclusion of projects/studies in the meta-analysis. However this approach is strongly criticized by the literature (Terrin 2003, Simonson et al 2014) whereby meta-analysts are not recommended to perform the trim-and-fill method when using meta-analysis software programs (Shi et al, 2019²⁰), as outliers and the pre-specified direction of missing studies could have influential impact on the trim-and-fill results. In addition a known limitation of Trim-and-Fill is that it can correct for publication bias that does not exist, underestimating effect sizes (Terrin et al 2003).

Although results adjusted using this approach remain largely positive they are presented in the Annex for the above mentioned considerations.

²⁰ Ref : need to put all references in footnote.

8 Implications for IFAD11

Turning to IFAD11 – what are the implications moving forward? In the following sections the process for selecting the IFAD11 IA is summarized, and descriptive analyses are presented, to assess for the presence of selection at the time of the projects' choice in the IFAD11 context, using performance ratings and other features of the portfolio universe.

Management has formalized the process of identifying candidate IFAD-supported projects to undergo ex post impact assessments. All regional divisions have been requested to identify and select potential countries and projects to conduct impact assessments from a list of all projects scheduled to close during IFAD11 (between 2019 and 2021) as of July 2018. Projects have been identified and selected through a participatory approach which involved Management and specifically the Research and Impact Assessment Division (RIA) and each of the five regional divisions, similar to the one implemented during IFAD10.

A first screening was done in July 2018 based on disbursement rate, timing of the project, and type of project. After this first screening, further identification was conducted based on learning potential, feasibility of conducting impact assessment given the eligibility and targeting criteria and project implementation, quality of M&E data, number of beneficiaries, type of interventions, and buy-in from country and project teams. RIA staff met with representatives with each regional division to select IFAD11 ex-post impact assessments.

During this meeting, each regional division received a list of projects that RIA staff had pre-screened²¹. RIA staff requested each regional division to identify six projects as candidates for impact assessments during IFAD11 (two projects per replenishment year).

Subsequently, a validation exercise was conducted through follow-up meetings in collaboration with each regional division and projects received clearance from both Country and Regional Directors. Additionally, RIA held internal discussions to ensure that projects selected were representative of the IFAD11 portfolio in terms of both regional distribution and sector.

8.1 IFAD11 Sample Validation

As part of the IFAD11 impact assessment agenda, 24 out of 121 projects have been selected for rigorous impact assessment equalling 19.8% of total projects, 20.7% of total financing, and 23.3% of total IFAD financing. Of the 121 projects belonging to the IFAD11 universe, nine²² were projects already part of evaluations initiated during IFAD9 and IFAD10 whose closing dates now fall during IFAD11. This gives a final universe of 112 projects eligible for evaluation in IFAD11. Considering the latter, the projects selected to be evaluated during IFAD11 account for 21.4% of the portfolio, representing 20.9% of total financing and 25.6% of IFAD financing.

²¹ The number of pre-screened projects scheduled to close between 2019 and 2021 that the RIA team had initially offered to each regional division were as follows: 26 projects for APR, 21 projects for ESA, 23 projects for LAC, 17 projects for NEN, and 16 projects for WCA.

²² The universe of "121 projects" include all projects CLOSING during IFAD11. However, upon further scrutiny it appeared that in the UNIVERSE of 121 – there are 9 projects whose evaluations were carried out during IFAD9 and IFAD10. The IFAD9 & IFAD10 projects are those whose closing dates were extended into IFAD11. The projects evaluated during IFAD10 are: Sao Tome and Principe PAPAC 1100001687; Senegal PAFA (extended as PAFA-E) 1100001693; Rwanda PRICE 1100001550; Kenya SDCP 1100001305; Nepal HVAP 1100001471; Bangladesh CCRIP 1100001647. The projects evaluated in IFAD9 then extended are : Uganda VODP2 1100001468; Ghana GASIP 1100001678; Bangladesh PACE 1100001648.

Table 7 presents the IFAD11 projects selected for impact assessment and their distribution by region, country, and project sector or type.

Table 7: IFAD11 Impact Assessments by Region, Country, Project Sector and Name

	<i>Region</i>	<i>country</i>	<i>sector</i>	<i>Project name</i>		<i>REGION</i>	<i>country</i>	<i>sector</i>	<i>Project name</i>
1	APR	India	CREDI	PT-Tamil Nadu	13	LAC	Argentina	MRKTG	PRODERI
2	APR	Pakistan	RURAL	SPPAP - PK	14	LAC	Bolivia	RURAL	ACCESOS
3	APR	Papua New Guinea	AGRIC	PPAP	15	LAC	Peru	RSRCH	PSSA
4	APR	Philippines	FISH	FishCORAL	16	NEN	Djibouti	RURAL	PRAREV-PECHE
5	APR	Solomon Islands	RURAL	RDP II	17	NEN	Kyrgyzstan	LIVST	LMDP
6	APR	Viet Nam	RURAL	AMD	18	NEN	Moldova, Republic of	RURAL	IRECR
7	ESA	Kenya	AGRIC	UTaNRMP	19	NEN	Morocco	AGRIC	PDFAZMH
8	ESA	Lesotho	RURAL	SADP	20	NEN	Tunisia	AGRIC	PRODESUD II
9	ESA	Malawi	RSRCH	SAPP	21	WCA	Ghana	CREDI	REP
10	ESA	Mozambique	AGRIC	PROSUL	22	WCA	Mali	CREDI	Rural Microfinance Programme
11	ESA	Tanzania, Un. Rep. of	MRKTG	MIVARF	23	WCA	Mauritania	RURAL	PASK II
12	ESA	Zambia	RSRCH	S3P	24	WCA	Nigeria	AGRIC	VCDP

Table 8 compares the regional distribution of these projects selected for impact assessment to the regional distribution of projects in the universe. Specifically, six projects were selected for impact assessment in both APR and ESA, five in NEN, four in WCA, and three in LAC. In the universe of 112 projects, there are 23 projects in LAC and NEN respectively, 32 in APR, 20 in ESA, and 17 in WCA.

Table 8: Distribution of projects in the universe and in the IA sample by Region

UNIVERSE BY REGION			IAS BY REGION		
REGION	Projects	%	REGION	Projects	%
APR	29	25.89	APR	6	25.00
ESA	20	17.86	ESA	6	25.00
LAC	23	20.54	LAC	3	12.50
NEN	23	20.54	NEN	5	20.83
WCA	17	15.18	WCA	4	16.67
TOTAL	112	100	Total	24	100.00

Table 9, first column, presents the current regional distribution of impact assessments in the IFAD11 sample and compares this distribution with three others, one weighted by the actual proportion of projects in each region and two more weighted by the actual proportional allocation of projects by financing and IFAD financing in the portfolio. As

shown in the table, the distribution of projects dictated by these different weighting schemes differ slightly from the distribution of those actually selected.

Table 9: Distribution of IFAD11 IAs. Current and proportional to project numbers and financing (total and IFAD only).

	Actual IFAD11 IA Distribution	Distribution by # of Projects	Distribution by Financing	Distribution by IFAD Financing
APR	6	6	8	7
ESA	6	4	6	5
LAC	3	5	3	3
NEN	5	5	3	4
WCA	4	4	4	5
	24	24	24	24

IFAD projects are also classified into eight project sectors (or types) in the universe (Table 10). IFAD11 projects are concentrated in rural development (46), agricultural development (29), and credit provision (16). The number of projects in the remaining five sectors range from two to six with the lowest concentration of projects in fisheries. The projects selected for IAs are comparatively less concentrated; eight rural development projects, six agricultural development projects, and three credit projects were selected. Despite the large number of credit projects in the portfolio, an equal number of projects (3) was selected in the research category along with two market access projects and one each in fisheries and livestock, respectively. No irrigation projects were selected for impact assessment during IFAD11.

Table 10: Distribution of projects in the universe and in the IA sample by Project Sector

Universe Sector	by Project		IAs Sector	sample by Project	
Sector	Projects	%	Sector	Projects	%
AGRIC	29	25.89	AGRIC	6	25
CREDI	16	14.29	CREDI	3	12.5
FISH	2	1.79	FISH	1	4.17
IRRIG	4	3.57	IRRIG	0	0
LIVST	6	5.36	LIVST	1	4.17
MRKTG	5	4.46	MRKTG	2	8.33
RSRCH	4	3.57	RSRCH	3	12.5
RURAL	46	41.07	RURAL	8	33.33
Total	112	100	Total	24	100

Looking at Table 11, note how the distribution of the IFAD11 IAs sample by sector differs from what would be dictated by the proportion of projects by sector in the universe, as well as the sectoral distribution proportional to the amount of total financing and IFAD financing. Certainly, when considering the projects by their sectoral classifications, irrigation projects seem to be underrepresented in the IFAD11 IA selection.

Table 11: Distribution of IFAD11 IAs by project sector. Current and proportional to project numbers and financing (total and IFAD only) by sector.

	IFAD11 IAs Sample Distribution	Distribution by # of Projects	Distribution by Financing	Distribution Total by Financing
AGRIC	6	6	6	7
CREDI	3	3	4	4
FISH	1	1	0	0
IRRIG	0	1	1	1
LIVST	1	1	1	1
MRKTG	2	1	1	1
RSRCH	3	1	2	1
RURAL	8	10	9	9
Total	24	24	24	24

Table 12 presents descriptive statistics of the sample and the unselected projects summarizing projects by their number of beneficiaries, financing, components, and implementation score. It also presents the results of t-tests for difference in means between the sample and unselected projects to assess for the possible presence of selection bias across the various samples.

Table 12: Balance test: selected indicators for IFAD11 (baseline).

	IAs Sample	Unselected Projects	Sample - Unselected	
	Mean	Mean	Diff. in Means	p-score
Beneficiaries	476,109	1,033,194	-557,085	0.595
Financing	81,820,640	74,454,692	7,365,947	0.695
IFAD Financing	37,848,004	29,311,323	8,536,680	0.176
# of Financiers	2.67	2.59	0.08	0.479
# of Subcomponent	5.92	6.10	-0.18	0.753

On average, the sample of IFAD11 IAs has 476,109 beneficiaries, \$81.8 million in financing, \$37.8 million in IFAD financing, 2.6 types of financiers, and 6 subcomponents. Compared to the average across the universe of IFAD11 projects, there are 942,862 beneficiaries, \$78.7 million in financing, \$31.9 million in IFAD financing, 2.6 types of financiers, and 6 subcomponents. Note how there are no statistically significant differences across the variables presented in Table 12, across the universe and the unselected projects.

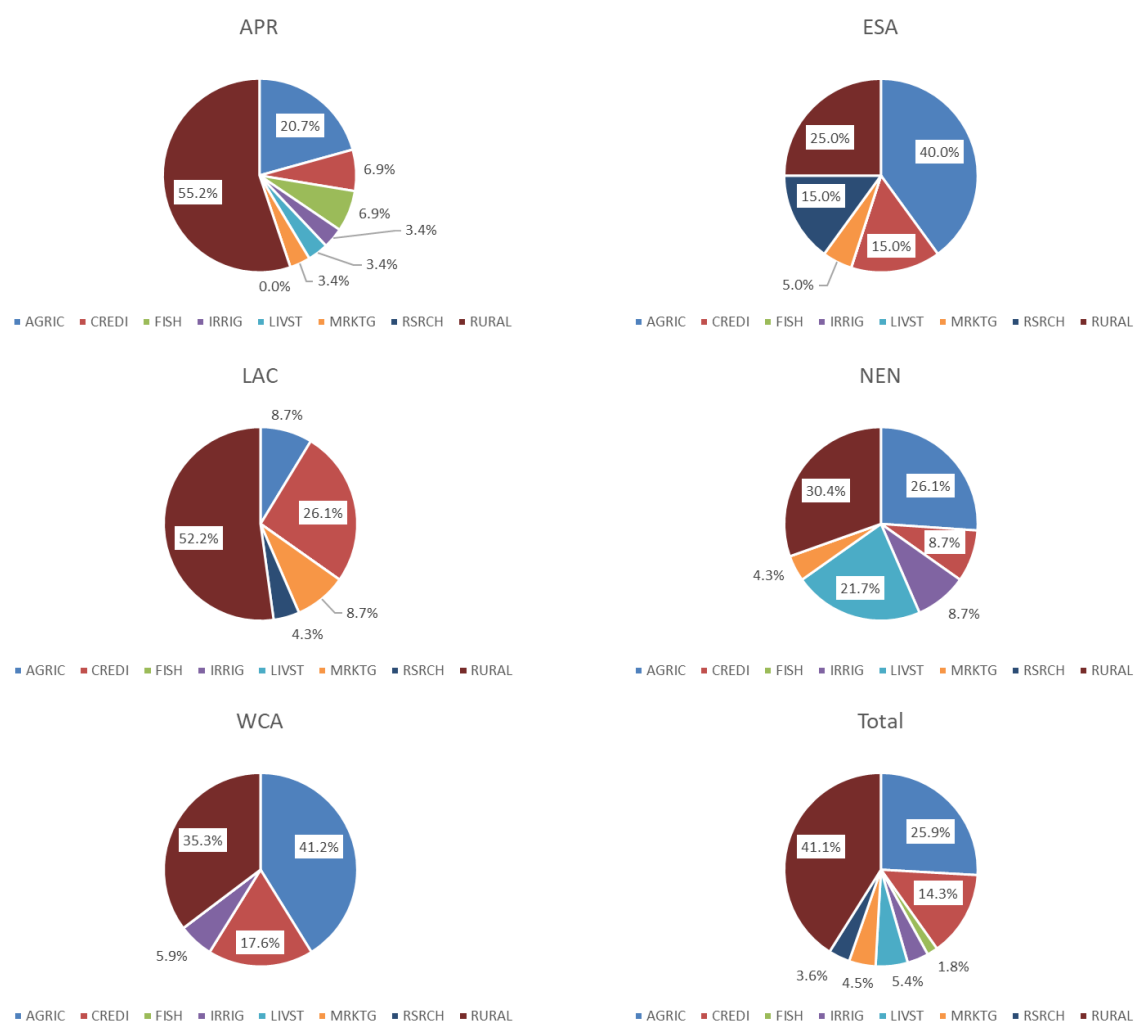
Table 13: Distribution of IFAD11 Universe by Project Type and Region

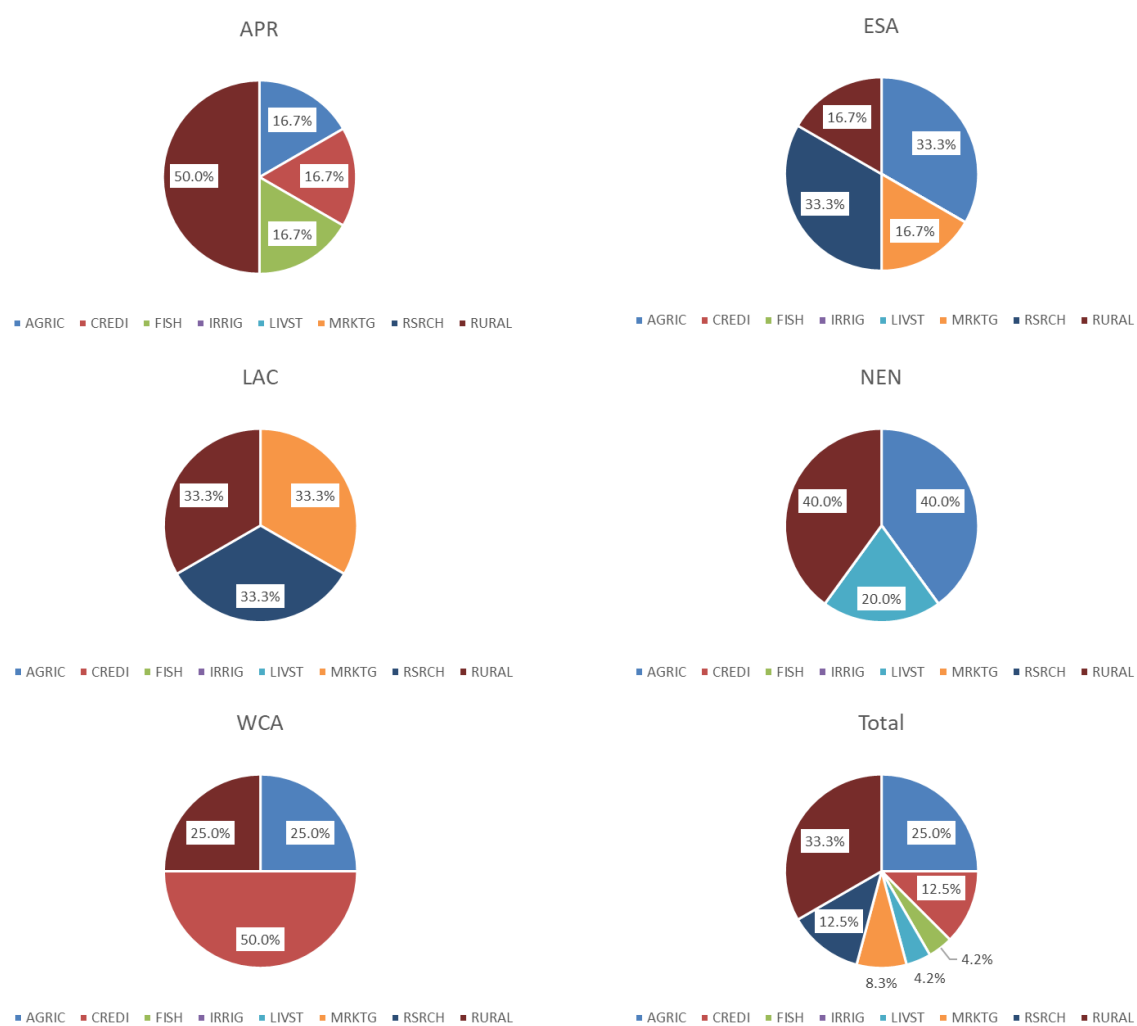
Table 14: Distribution of IFAD11 IA Sample by Project Type and Region

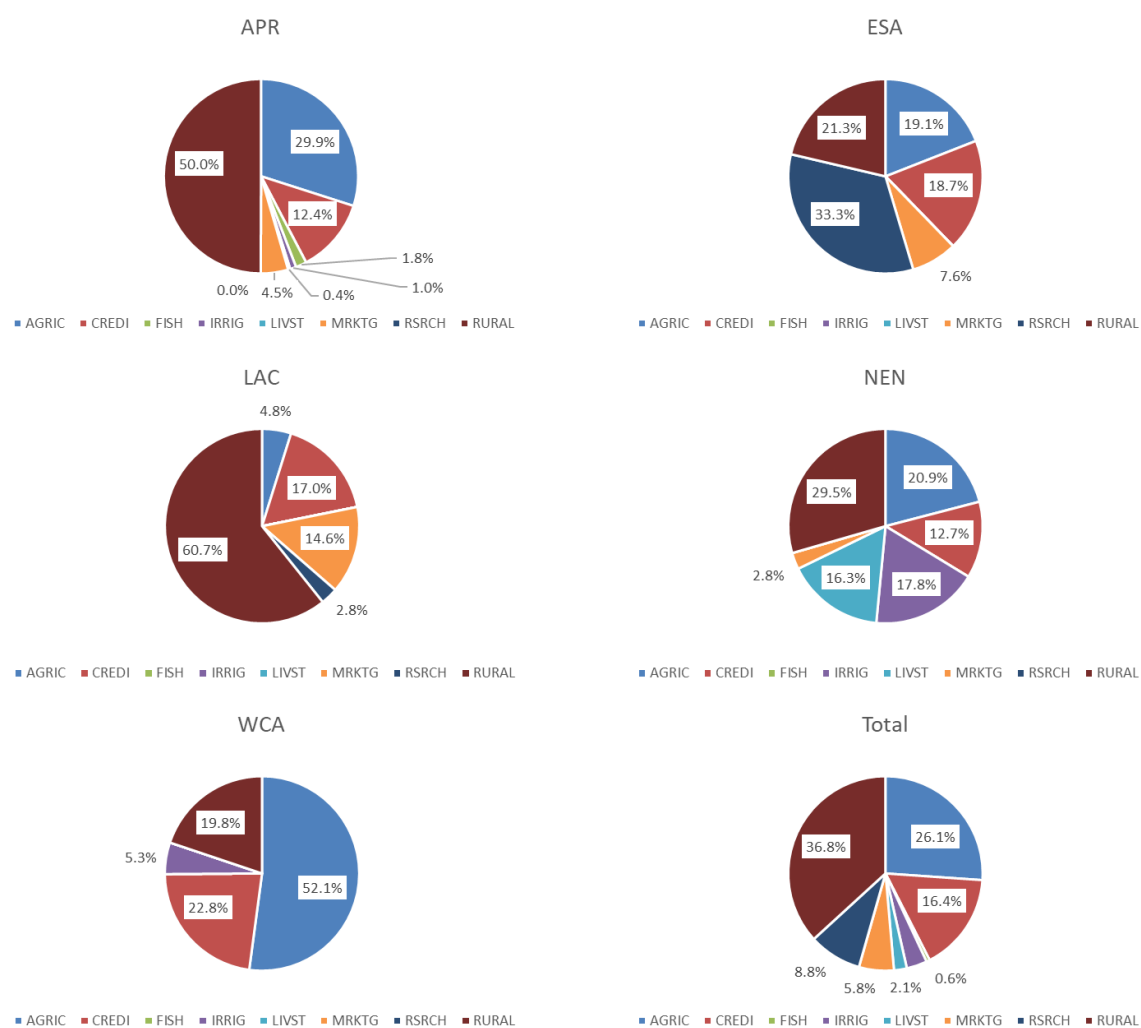
Table 15: Distribution of IFAD11 Universe by Project Type and Region Financing

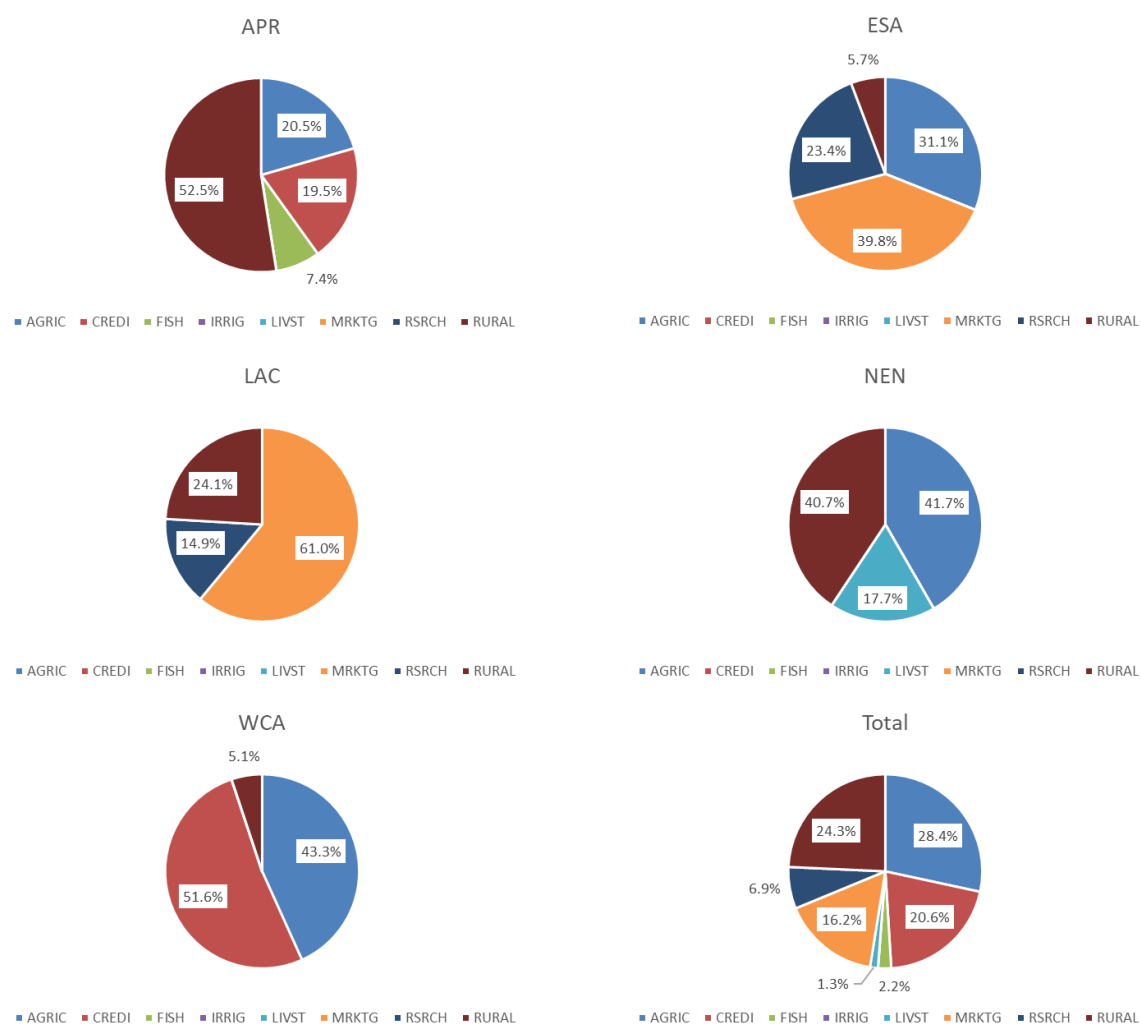
Table 16: Distribution of IFAD11 IA Sample by Project Type and Region Financing

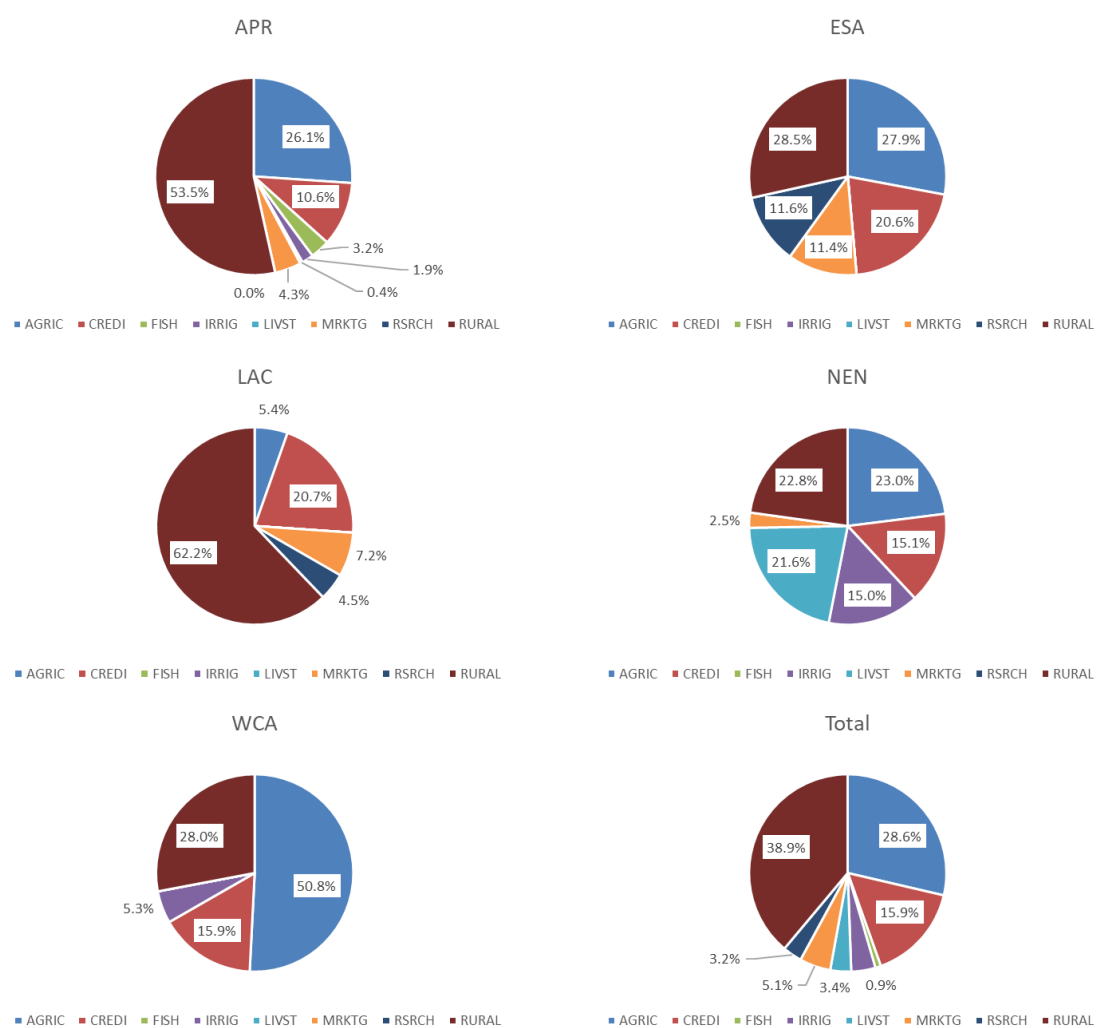
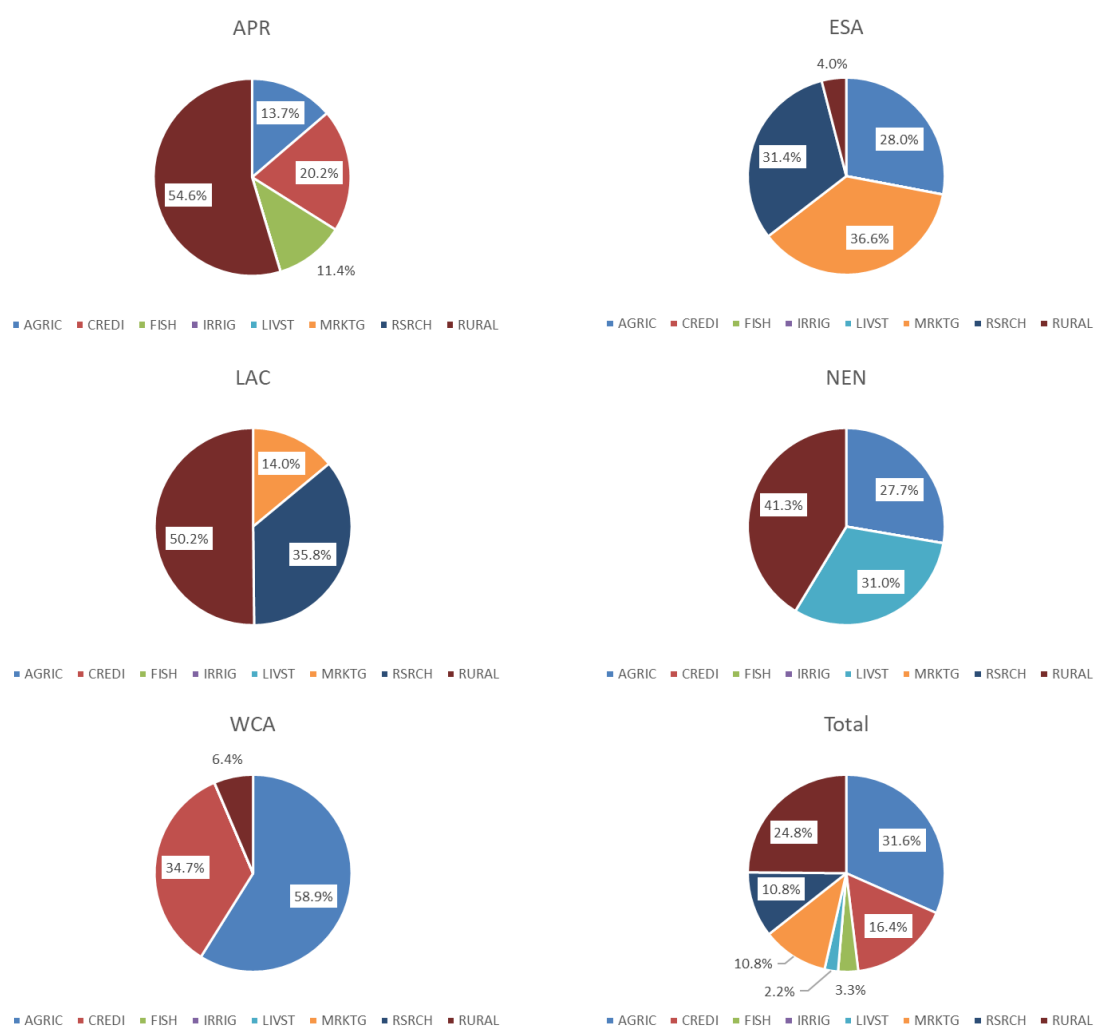
Table 17: Distribution of IFAD11 Universe by Project Type and Region IFAD Financing

Table 18: Distribution of IFAD11 IA Sample by Project Type and Region IFAD Financing

In addition to being divided into project sector and regional categories, projects can further be disaggregated into component, subcomponent and subcomponent type. The only variable that is standardized in the system (GRIPS) is the subcomponent type, which can be broken down into approximately 60 categories. Note that the component feature is not standardized across IFAD databases and therefore cannot be used in this analysis.

Each project contains multiple subcomponents funded by different sources. Nearly all projects contain a management and coordination component, but there remains considerable heterogeneity. Apart from the management subcomponent, the most common subcomponent in the universe of projects is rural financial services (44.7%), followed by local capacity building (38.39%), monitoring and evaluation (28.57%), and institutional support (27.68%). Table 19 compares the relative distribution of subcomponent types in the selected IAs and the IFAD11 portfolio.

When projects are categorized by their largest subcomponent (in terms of current financing), the most common subcomponents among IFAD11 IAs are development funds and rural financial services each with three projects selected for IAs, respectively (table available upon request). This matches the portfolio overall where development funds and rural financial services are the largest subcomponents with thirteen and fourteen projects, respectively. However, beyond that similarity, the distribution of IAs unevenly represents the distribution in the portfolio. Finally, looking at Table 20– the distribution of projects by largest financier and region is provided in the portfolio and in the IA sample.

Two issues in this analysis hinder both our ability to draw a representative sample by sector and components distribution and assess the representativeness of our current selection. The first is that the sector variable does not reflect the true nature of the project. For example, projects in which the true intervention is related to livestock or animal husbandry may be classified as marketing if there is substantial intervention in business formation or a rural development project may have a substantive irrigation component but not be classified as such because it is combined with other interventions. Moreover, the sectors are too broad such that there is substantial heterogeneity within sectors, namely the rural and agricultural development. Secondly, the subcomponent type is too disaggregated (60 unique entries over 112 records) to be used as a possible stratification feature hence it would need to be recoded prior to be used. However, there currently exists no method for standardizing or harmonizing subcomponents within projects and each project can have multiple subcomponents all of varying sizes and relative importance.

Table 19: Distribution of subcomponent types in the IAs sample and in the IFAD11 universe

Subcomponent Type	IAs		Universe	
	Freq.	%	Freq.	%
<i>Rural financial services</i>	10	41.67	50	44.64
<i>Local capacity building</i>	8	33.33	43	38.39
<i>Monitoring and evaluation</i>	7	29.17	32	28.57
<i>Institutional support</i>	5	20.83	31	27.68
<i>Technology transfer</i>	7	29.17	27	24.11
<i>Rural infrastructure</i>	5	20.83	25	22.32
<i>Business development</i>	2	8.33	21	18.75
<i>Development funds</i>	5	20.83	20	17.86
<i>Irrigation infrastructure</i>	1	4.17	19	16.96
<i>Community development</i>	5	20.83	18	16.07
<i>Crop production</i>	3	12.50	17	15.18
<i>Marketing: inputs/outputs</i>	6	25.00	17	15.18

<i>Roads/tracks</i>	1	4.17	17	15.18
<i>Micro-enterprises</i>	3	12.50	16	14.29
<i>Climate change adaptation</i>	4	16.67	14	12.50
<i>Rangelands/pastures</i>	2	8.33	14	12.50
<i>Credit</i>	1	4.17	13	11.61
<i>Resource mgmt/protection</i>	1	4.17	13	11.61
<i>Rural enterprises</i>	2	8.33	13	11.61
<i>Animal husbandry</i>	2	8.33	11	9.82
<i>Technology development</i>	2	8.33	11	9.82
<i>Policy Support/Development</i>	2	8.33	9	8.04
<i>Animal health</i>	1	4.17	8	7.14
<i>Crop extension services</i>	2	8.33	8	7.14
<i>Market information/study</i>	3	12.50	8	7.14
<i>Training</i>	3	12.50	8	7.14
<i>Communication</i>	1	4.17	7	6.25
<i>Market infrastructure</i>	3	12.50	7	6.25
<i>Soil and Water conservation</i>	1	4.17	7	6.25
<i>Livestock - other</i>	0	0.00	6	5.36
<i>Livestock post-harvest</i>	2	8.33	6	5.36
<i>Seed, fertilizer, pesticide</i>	1	4.17	6	5.36
<i>Health and nutrition</i>	0	0.00	5	4.46
<i>Drinking water/sanitation</i>	0	0.00	4	3.57
<i>Forestry</i>	1	4.17	4	3.57
<i>Literacy</i>	0	0.00	4	3.57
<i>Crop technology development</i>	1	4.17	3	2.68
<i>Disaster mitigation</i>	2	8.33	3	2.68
<i>Financing/preparation charges</i>	1	4.17	3	2.68
<i>Fisheries infrastructure</i>	2	8.33	3	2.68
<i>Irrigation management</i>	0	0.00	3	2.68
<i>Land improvement</i>	1	4.17	3	2.68
<i>On-farm storage</i>	1	4.17	3	2.68
<i>Animal restocking</i>	1	4.17	2	1.79
<i>Aquaculture</i>	0	0.00	2	1.79
<i>Education (primary/second)</i>	1	4.17	2	1.79
<i>Fisheries/marine conservation</i>	2	8.33	2	1.79
<i>Input supply</i>	1	4.17	2	1.79
<i>Knowledge management</i>	0	0.00	2	1.79
<i>Land reform/titles</i>	0	0.00	2	1.79
<i>Processing</i>	0	0.00	2	1.79
<i>Standards and regulations</i>	1	4.17	2	1.79
<i>Energy production</i>	0	0.00	1	0.89
<i>Fishing (capture)</i>	0	0.00	1	0.89
<i>Housing</i>	1	4.17	1	0.89
<i>Insurance/risk transfer</i>	1	4.17	1	0.89
<i>Legal assistance</i>	0	0.00	1	0.89
<i>Mechanization services</i>	0	0.00	1	0.89
<i>Venture capital</i>	1	4.17	1	0.89

Finally, looking at Table 20 – the distribution of projects by largest financier and region is provided in the portfolio and in the IA sample.

Table 20 Distribution of Largest Financier Type by Region.

IFAD11 Universe

	APR	ESA	LAC	NEN	WCA	Total
<i>Largest Financier</i>						
<i>Domestic</i>	10	4	9	3	1	27
<i>IFAD</i>	13	13	12	18	14	70
<i>International</i>	6	3	2	2	2	15
<i>Total</i>	29	20	23	23	17	112

IFAD11 IAs Sample

	APR	ESA	LAC	NEN	WCA	Total
<i>Largest Financier</i>						
<i>Domestic</i>	1	0	1	2	1	5
<i>IFAD</i>	3	5	2	3	3	16
<i>International</i>	2	1	0	0	0	3
<i>Total</i>	6	6	3	5	4	24

Table **21** presents additional baseline ratings statistics based on the first available implementation performance rating available in the system - given by the corresponding supervision report. Here, only performance of M&E system seems to be statistically significant. However, given that 24 implementation ratings have been tested - these results stress the absence of selection bias in the case of the IFAD11 sample of IAs.

Table 21: Balance tests: implementation performance ratings for IFAD11 (baseline characteristics)

	IFAD11 IA Sample	Unselected Projects (closing during IFAD11)	IFAD 11 Universe	Sample - Unselected		Sample - Universe	
	Mean	Mean	Mean	Diff. in Means	p- score	Diff. in Means	p- score
<i>Assessment of the Overall Implementation Performance</i>	3.96	3.94	3.94	0.02	0.80	0.02	0.80
<i>Likelihood of Achieving the Development Objective Effectiveness</i>	4.04	4.02	4.04	0.02	0.82	0.01	0.94
<i>Targeting and Outreach</i>	3.91	3.96	3.95	-0.05	0.61	-0.04	0.69
<i>Gender equality & women's participation</i>	4.08	4.09	4.09	0.00	0.98	-0.01	0.95
<i>Agricultural Productivity</i>	4.00	4.04	4.04	-0.04	0.71	-0.04	0.62
<i>Adaptation to Climate Change</i>	4.00	3.97	3.98	0.03	0.66	0.02	0.78
<i>Institutions and Policy Engagement</i>	4.08	4.00	4.02	0.08	0.40	0.07	0.48
<i>Human and Social Capital and Empowerment</i>	4.00	3.97	3.99	0.03	0.73	0.01	0.92
<i>Quality of Beneficiary Participation</i>	3.95	4.03	4.02	-0.07	0.39	-0.07	0.26
<i>Responsiveness of Service Providers</i>	4.04	4.05	4.05	-0.01	0.94	-0.01	0.89
<i>Environment and Natural Resource Management</i>	4.00	3.96	3.97	0.04	0.68	0.03	0.70
<i>Exit Strategy</i>	4.09	3.98	4.00	0.11	0.21	0.09	0.37
<i>Potential for Scaling-up</i>	3.93	4.00	3.99	-0.07	0.40	-0.05	0.48
<i>Quality of Project Management</i>	3.95	4.07	4.04	-0.12	0.15	-0.09	0.13
<i>Knowledge Management</i>	3.96	3.96	3.97	-0.01	0.97	-0.02	0.89
<i>Coherence between AWPB and Implementation</i>	4.00	3.99	3.99	0.01	0.87	0.01	0.89
<i>Performance of M&E System</i>	3.92	3.77	3.82	0.15	0.24	0.09	0.35
<i>Acceptable Disbursement Rate</i>	4.04	3.73	3.80	0.31	0.00	0.24	0.01
<i>Quality of Financial Management</i>	2.71	3.10	3.02	-0.39	0.31	-0.31	0.44
<i>Quality and Timeliness of Audit</i>	3.96	4.00	3.98	-0.04	0.71	-0.02	0.79
<i>Counterparts Funds</i>	4.00	3.93	3.95	0.07	0.38	0.05	0.44
<i>Compliance with Loan Covenants</i>	4.08	3.99	4.01	0.10	0.50	0.07	0.64
<i>Procurement</i>	4.00	3.98	3.99	0.02	0.83	0.01	0.94
	3.83	3.94	3.92	-0.11	0.33	-0.09	0.53

8.2 Conclusions for IFAD11

In light of these descriptive analyses and broad considerations, Management can conclude that there is no selection bias in the IFAD11 sample of impact assessments.

However the additional following recommendation can be made, notably adjusting the regional distribution to allow for two more impact assessments in LAC.

Appendix – Annex I

Table 22: Balance test: Implementation Performance Ratings (Baseline Characteristics) by region

	APR						ESA					
	Sample		Unselected		Sample - Unselected		Sample		Unselected		Sample - Unselected	
	n	Mean	n	Mean	Difference	p-score	n	Mean	n	Mean	Difference	p-score
Duration	5	6.80	25	8.56	-1.76	0.31	6	10.00	14	8.29	1.71	0.07
Beneficiaries	5	1,226,531	25	1,108,359	118,173	0.92	6	377,717	14	464,178	-86,461	0.85
Approved Funding	5	74,052,700	25	62,145,347	11,907,353	0.64	6	68,786,299	14	83,084,966	-83,084,966	0.85
Assessment of the Overall Implementation Performance	5	4.20	25	3.88	0.32	0.07	6	4.17	14	3.71	0.45	0.17
Likelihood of Achieving the Development Objective	5	4.00	25	3.92	0.08	0.72	6	4.00	14	3.86	0.14	0.61
Effectiveness	4	3.50	19	3.84	-0.34	0.14	3	4.00	13	3.92	0.08	0.65
Targeting and Outreach	5	4.20	25	3.96	0.24	0.19	6	4.17	14	4.00	0.17	0.40
Gender equality & women's participation	5	4.20	25	4.00	0.20	0.33	6	4.17	14	4.07	0.10	0.54
Agricultural Productivity	5	4.40	19	4.00	0.40	0.00	5	4.00	14	3.71	0.29	0.32
Adaptation to Climate Change	1	4.00	3	4.00	0.00	-	-	-	-	-	-	-
Institutions and Policy Engagement	5	4.40	21	4.00	0.40	0.10	6	4.33	14	4.00	0.33	0.36
Human and Social Capital and Empowerment	4	4.25	21	3.95	0.30	0.36	6	4.17	14	3.79	0.38	0.16
Quality of Beneficiary Participation	5	4.00	25	4.04	-0.04	0.80	6	4.00	14	4.00	0.00	-
Responsiveness of Service Providers	5	4.40	25	3.84	0.56	0.01	6	3.67	14	4.07	-0.41	0.11
Environment and Natural Resource Management	1	4.00	3	4.00	0.00	-	-	-	-	-	-	-
Exit Strategy	4	4.25	19	4.00	0.25	0.03	3	4.00	3	3.67	0.33	0.37
Potential for Scaling-up	4	4.25	19	4.00	0.25	0.03	5	4.40	14	4.14	0.26	0.51
Quality of Project Management	5	4.00	25	3.84	0.16	0.61	6	4.17	14	3.64	0.52	0.24
Knowledge Management	4	4.00	22	3.86	0.14	0.45	6	4.17	13	3.92	0.24	0.14
Coherence between AWPB and Implementation	5	4.00	22	3.82	0.18	0.32	6	4.17	14	3.86	0.31	0.37
Performance of M&E System	5	3.80	24	3.88	-0.08	0.80	6	3.50	14	3.71	-0.21	0.47
Acceptable Disbursement Rate	5	4.20	25	3.32	0.88	0.14	6	4.17	14	3.86	0.31	0.62
Quality of Financial Management	5	4.40	22	3.91	0.49	0.04	5	3.80	14	3.71	0.09	0.81

Quality and Timeliness of Audit	5	4.40	24	3.92	0.48	0.07	6	4.00	14	3.86	0.14	0.53
Counterparts Funds	5	4.40	25	4.00	0.40	0.28	6	4.17	14	4.29	-0.12	0.80
Compliance with Loan Covenants	5	4.40	25	3.92	0.48	0.01	6	4.33	14	3.79	0.55	0.10
Procurement	5	4.40	25	3.96	0.44	0.03	6	3.83	14	3.79	0.05	0.88

	LAC						NEN					
	Sample		Unselected		Sample - Unselected		Sample		Unselected		Sample - Unselected	
	n	Mean	n	Mean	Difference	p-score	n	Mean	n	Mean	Difference	p-score
Duration	3	6.667	15	8.00	-1.33	0.35	1	7.00	12	7.92	-0.92	0.66
Beneficiaries	3	40,518	15	63,093	-22,576	0.62	1	145,600	12	92,023	53,577	0.58
Approved Funding	3	31,437,826	15	25,628,214	5,809,612	0.57	1	15,780,852	12	38,844,717	-23,063,865	0.31
Assessment of the Overall Implementation Performance	3	4	15	3.53	0.47	0.36	1	4.00	12	4.00	0.00	1.00
Likelihood of Achieving the Development Objective	3	4.00	15	3.73	0.27	0.46	1	4.00	12	4.17	-0.17	0.79
Effectiveness	1	4.00	10	3.60	0.40	0.66	1	4.00	8	3.75	0.25	0.75
Targeting and Outreach	3	4.67	15	4.13	0.53	0.20	1	4.00	12	4.42	-0.42	0.45
Gender equality & women's participation	3	3.67	15	3.67	0.00	1.00	1	4.00	12	4.08	-0.08	0.88
Agricultural Productivity	1	4.00	10	3.90	0.10	0.90	1	4.00	10	3.90	0.10	0.87
Adaptation to <i>Climate Change</i>	-	-	-	-	-	-	-	-	-	-	-	-
Institutions and Policy Engagement	3	3.33	11	3.82	-0.49	0.44	1	4.00	12	3.92	0.08	0.88
Human and Social Capital and Empowerment	1	4.00	11	3.91	0.09	0.92	1	4.00	12	3.92	0.08	0.88
Quality of Beneficiary Participation	3	3.67	15	3.87	-0.20	0.74	1	4.00	12	4.08	-0.08	0.88
Responsiveness of Service Providers	3	3.67	15	3.73	-0.07	0.83	1	4.00	12	4.08	-0.08	0.91
Environment and Natural Resource Management	-	-	-	-	-	-	-	-	-	-	-	-
Exit Strategy	1	4.00	9	3.67	0.33	0.67	1	4.00	11	4.09	-0.09	0.90
Potential for Scaling-up	1	4.00	10	3.80	0.20	0.81	1	4.00	11	4.27	-0.27	0.75
Quality of Project Management	3	4.00	15	3.67	0.33	0.60	1	3.00	12	4.08	-1.08	0.15
Knowledge Management	2	5.00	10	4.00	1.00	0.18	1	4.00	12	4.25	-0.25	0.71
Coherence between AWPB and Implementation	2	4.50	12	3.50	1.00	0.12	1	4.00	12	3.83	0.17	0.79

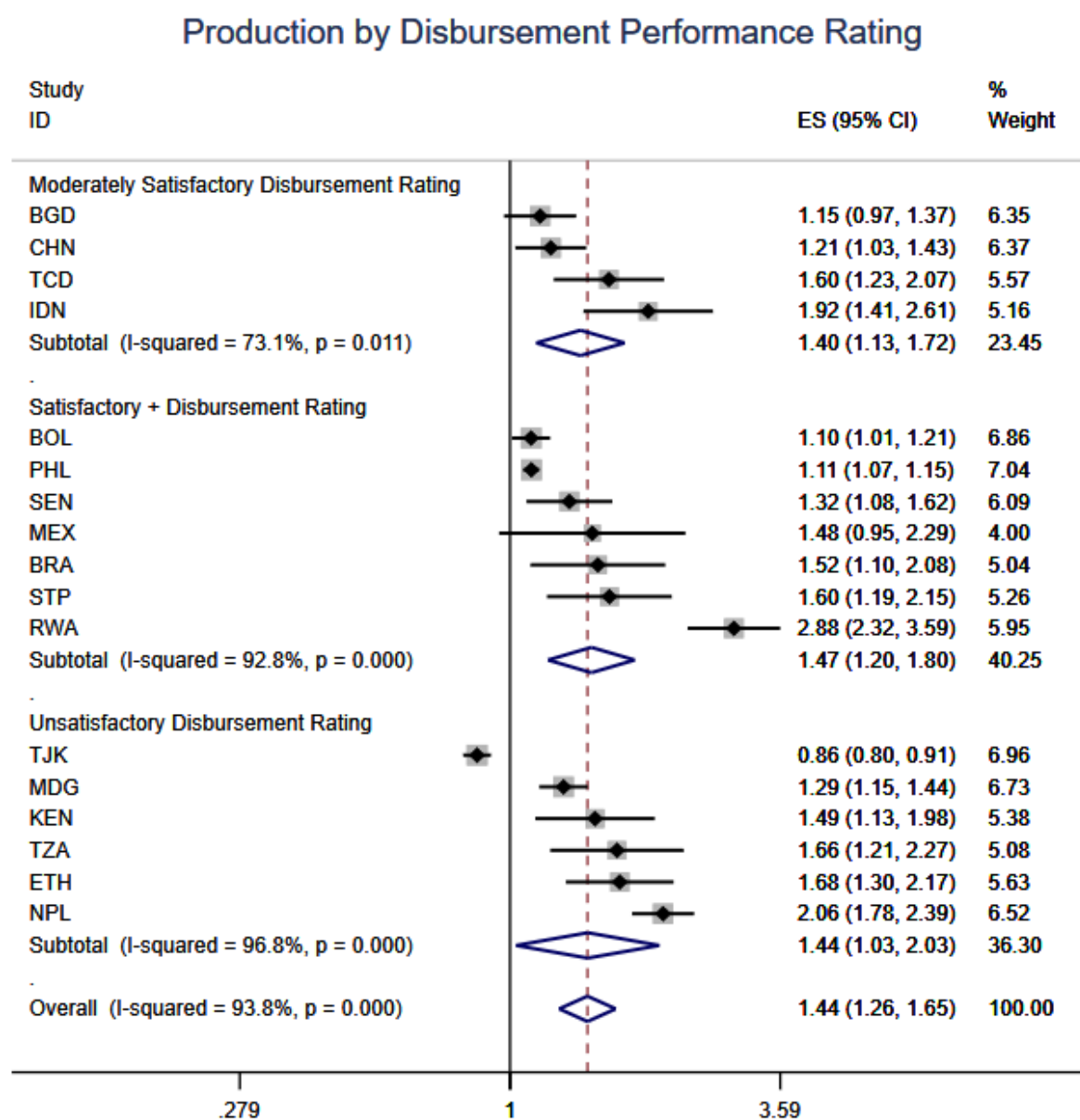
Performance of M&E System	3	3.67	15	3.53	0.13	0.78	1	4.00	12	4.08	-0.08	0.88
Acceptable Disbursement Rate	3	4.33	15	3.47	0.87	0.35	1	1.00	12	3.83	-2.83	0.14
Quality of Financial Management	3	4.00	13	3.85	0.15	0.83	1	4.00	11	4.18	-0.18	0.68
Quality and Timeliness of Audit	3	4.00	15	4.00	0.00	1.00	1	4.00	12	4.25	-0.25	0.76
Counterparts Funds	3	4.67	15	3.67	1.00	0.06	1	4.00	12	4.33	-0.33	0.63
Compliance with Loan Covenants	3	4.00	15	4.00	0.00	1.00	1	3.00	12	4.17	-1.17	0.02
Procurement	3	4.67	15	4.07	0.60	0.28	1	4.00	12	4.08	-0.08	0.92

	WCA											
	Sample		Unselected		Sample - Unselected							
	n	Mean	n	Mean	Difference	p-score						
Duration	4	8.50	22	8.32	0.18	0.90						
Beneficiaries	4	162583	22	897,950	-735,368	0.55						
Approved Funding	4	22364480	22	40,751,001	-18,386,521	0.34						
Assessment of the Overall Implementation Performance	4	4.00	22	4.05	-0.05	0.68						
Likelihood of Achieving the Development Objective	4	4.00	22	4.18	-0.18	0.37						
Effectiveness	3	4.00	19	4.11	-0.11	0.58						
Targeting and Outreach	4	4.25	22	4.18	0.07	0.80						
Gender equality & women's participation	4	4.00	22	4.09	-0.09	0.77						
Agricultural Productivity	3	4.00	18	4.11	-0.11	0.57						
Adaptation to Climate Change	1	4.00	3	4.33	-0.33	0.67						
Institutions and Policy Engagement	3	4.00	20	4.20	-0.20	0.52						
Human and Social Capital and Empowerment	3	4.00	19	4.00	0.00	1.00						
Quality of Beneficiary Participation	4	4.00	22	4.23	-0.23	0.41						
Responsiveness of Service Providers	4	3.75	22	4.14	-0.39	0.21						
Environment and Natural Resource Management	1	4.00	5	4.00	0.00	-						
Exit Strategy	2	4.00	16	4.06	-0.06	0.74						
Potential for Scaling-up	3	4.00	18	4.11	-0.11	0.57						
Quality of Project Management	4	3.75	22	4.00	-0.25	0.56						
Knowledge Management	3	4.00	19	4.16	-0.16	0.48						

Coherence between AWPB and Implementation	3	4.00	21	3.86	0.14	0.79						
Performance of M&E System	4	4.00	22	3.91	0.09	0.81						
Acceptable Disbursement Rate	4	5.00	22	3.05	1.96	0.02						
Quality of Financial Management	3	4.33	19	3.90	0.44	0.23						
Quality and Timeliness of Audit	4	4.00	22	4.09	-0.09	0.86						
Counterparts Funds	4	4.75	22	3.91	0.84	0.01						
Compliance with Loan Covenants	4	4.25	22	4.23	0.02	0.93						
Procurement	4	4.00	22	4.09	-0.09	0.68						

Appendix - Annex II

Table 23



Note: Based on inverse variance weighting

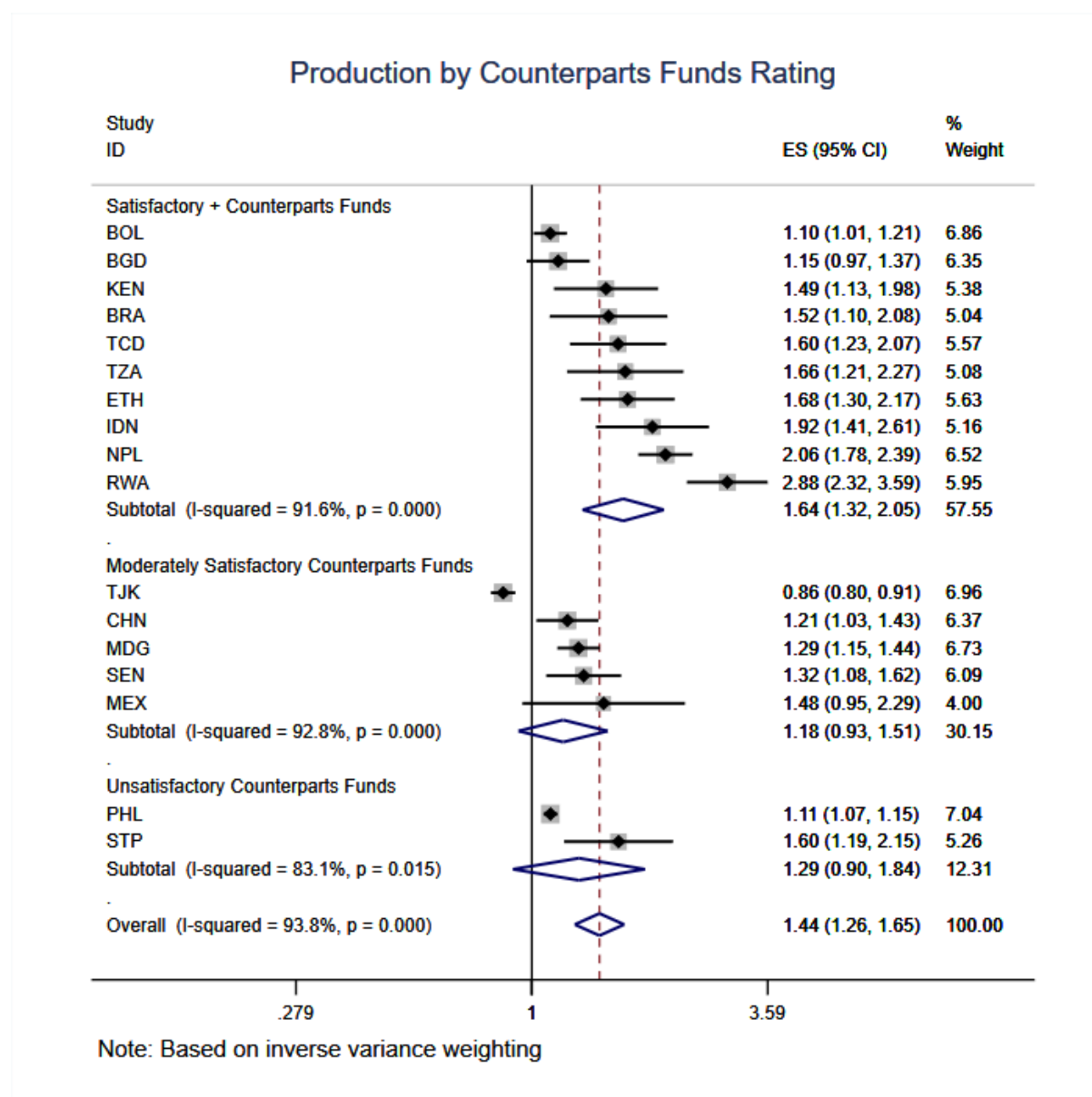
Table 24

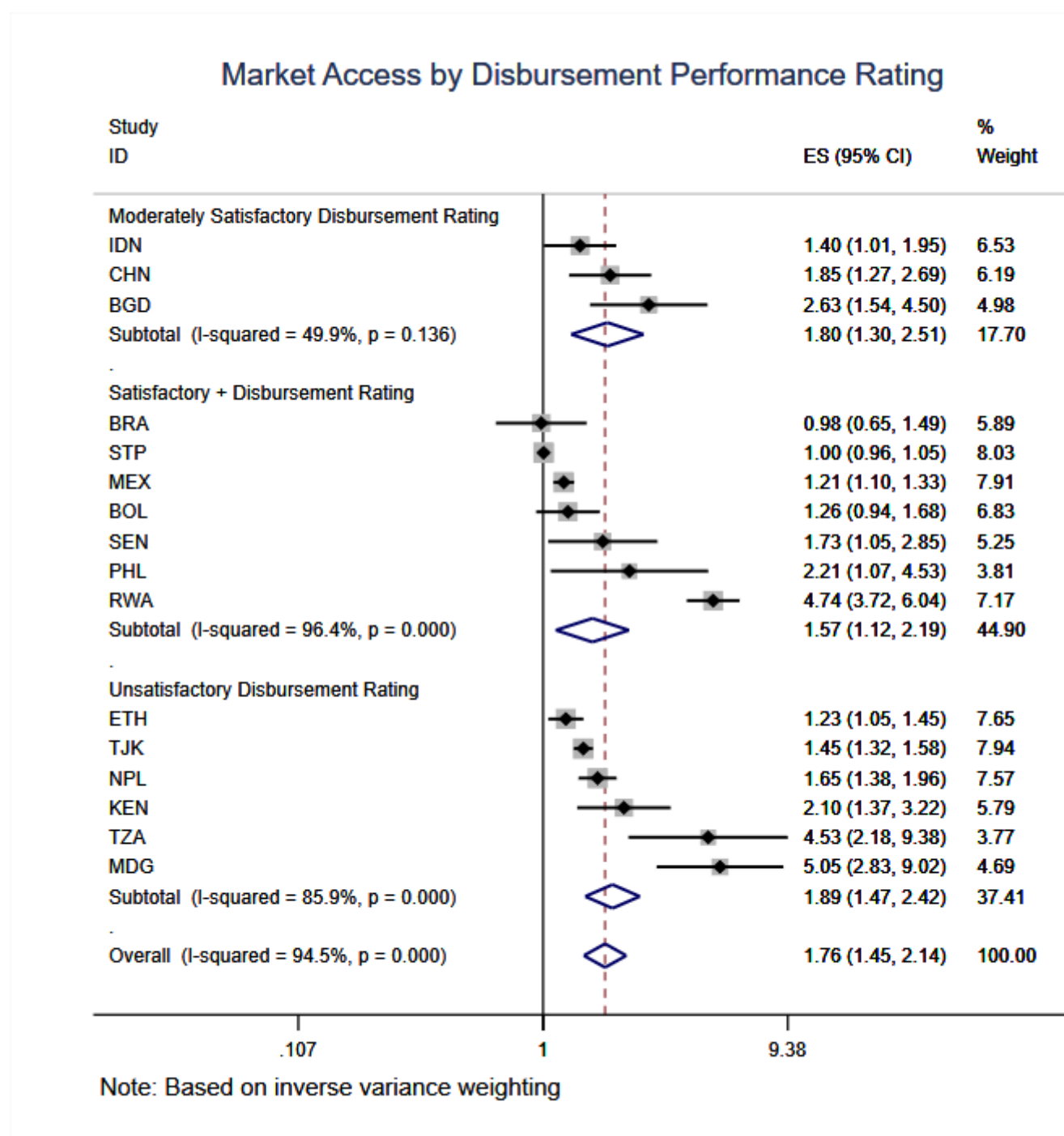
Table 25

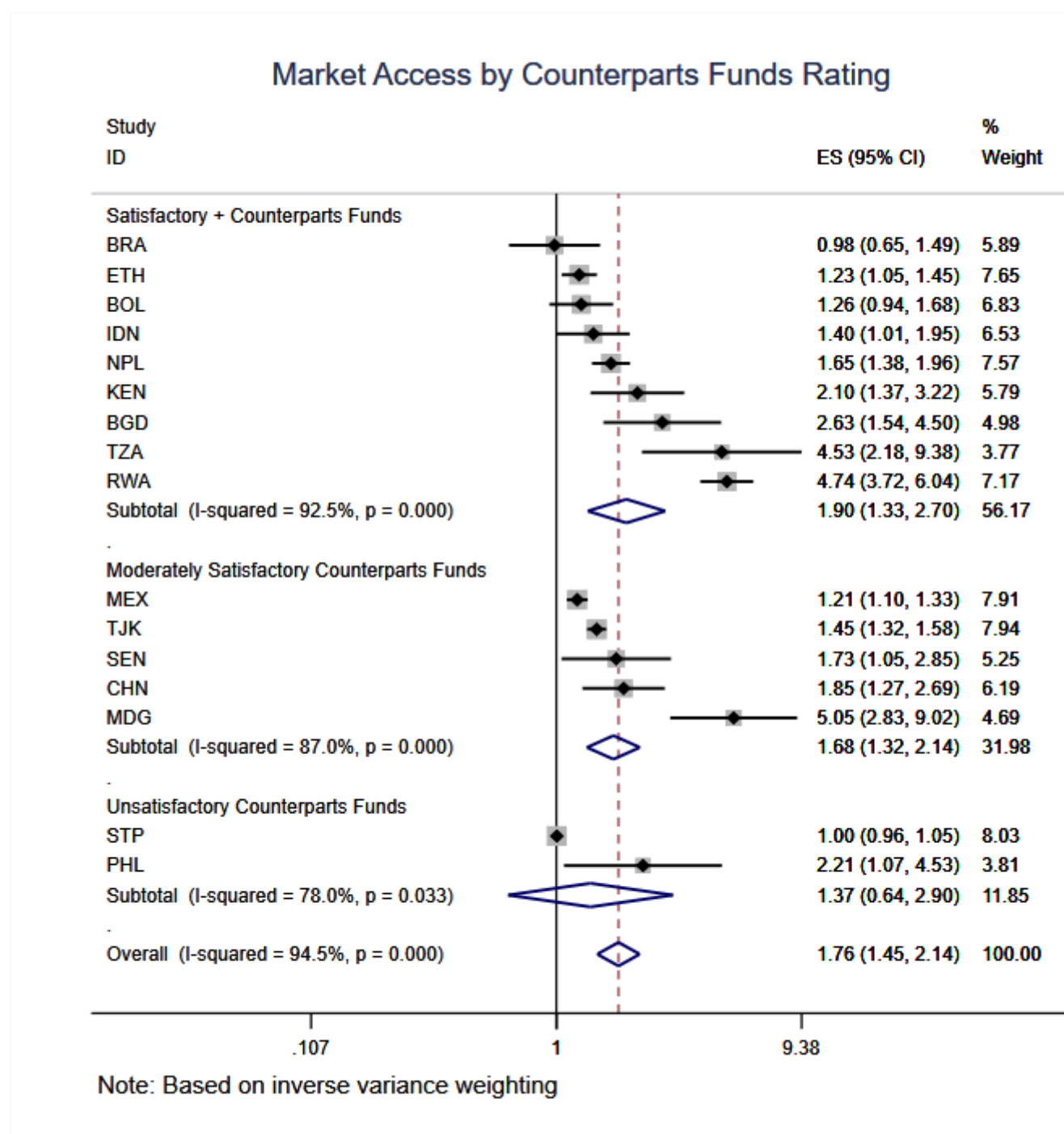
Table 26

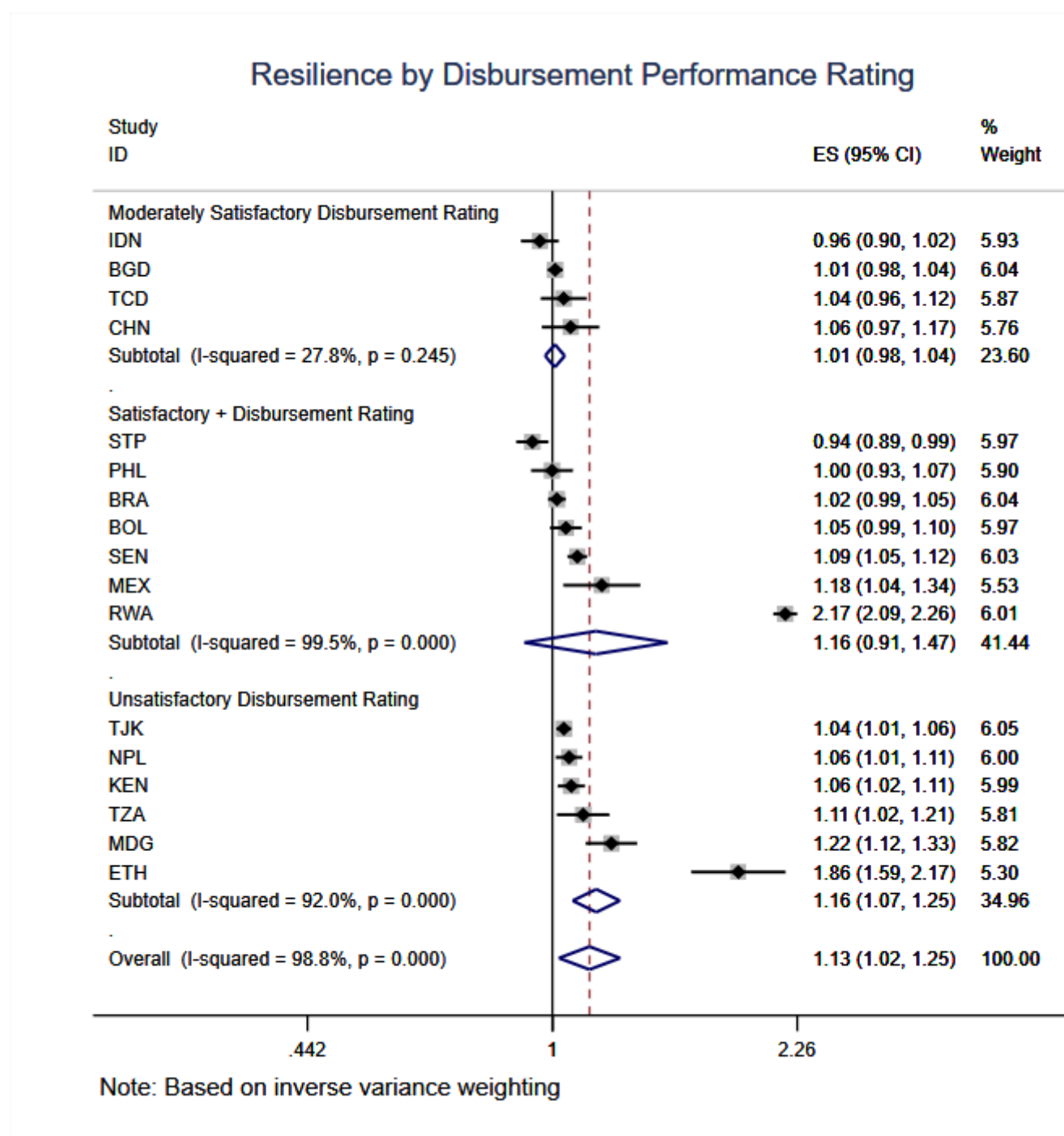
Table 27

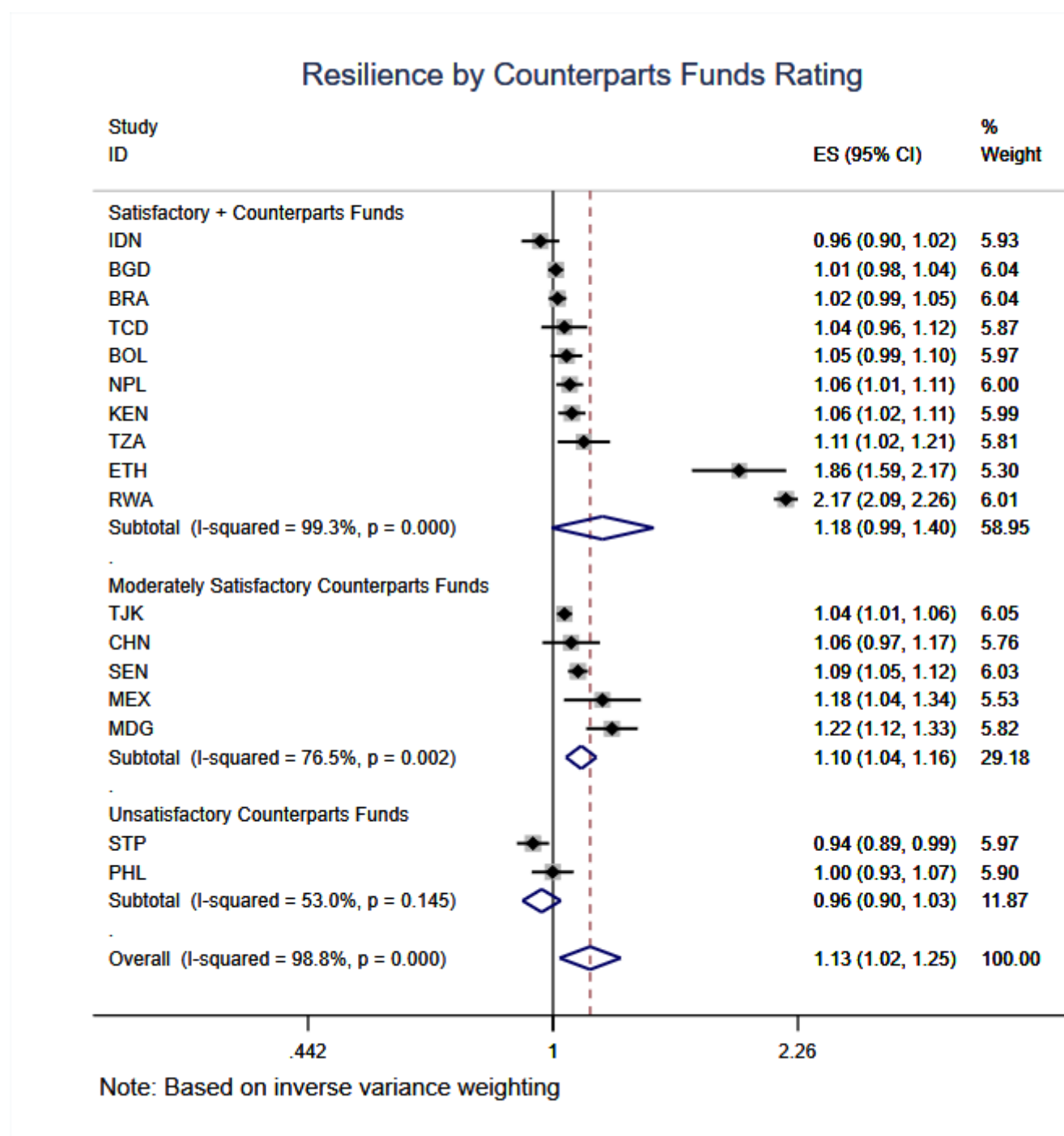
Table 28

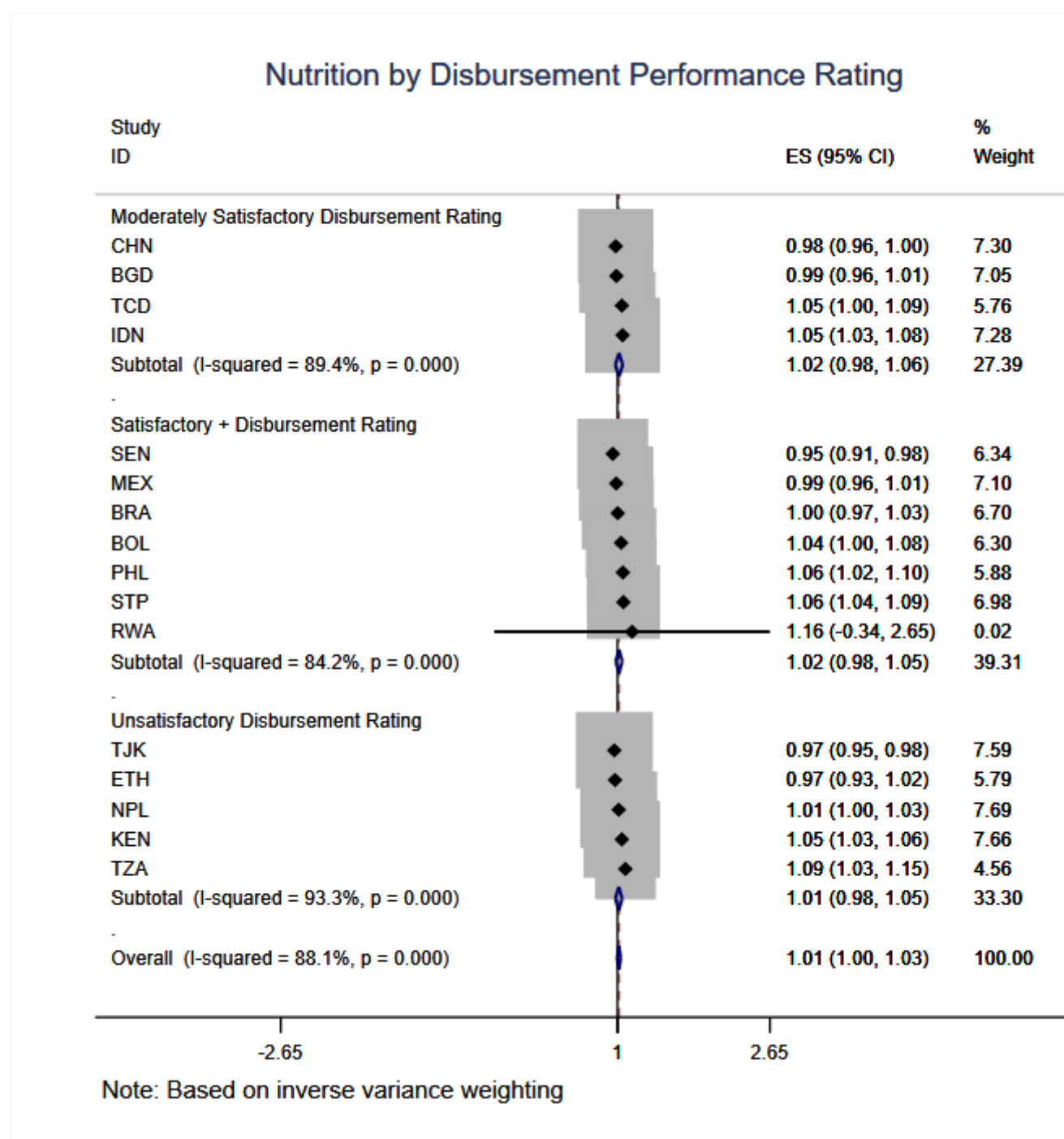
Table 29

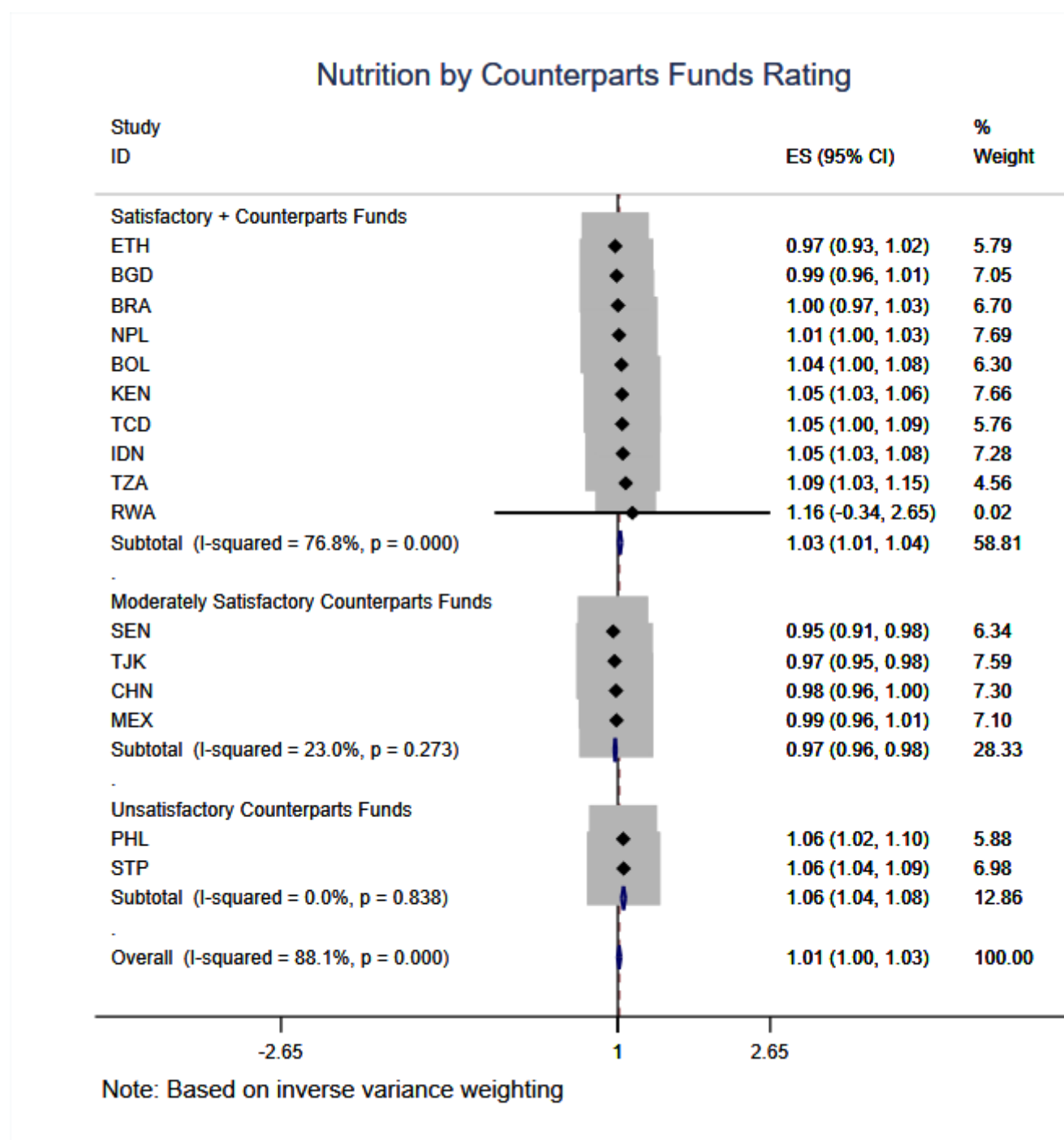
Table 30

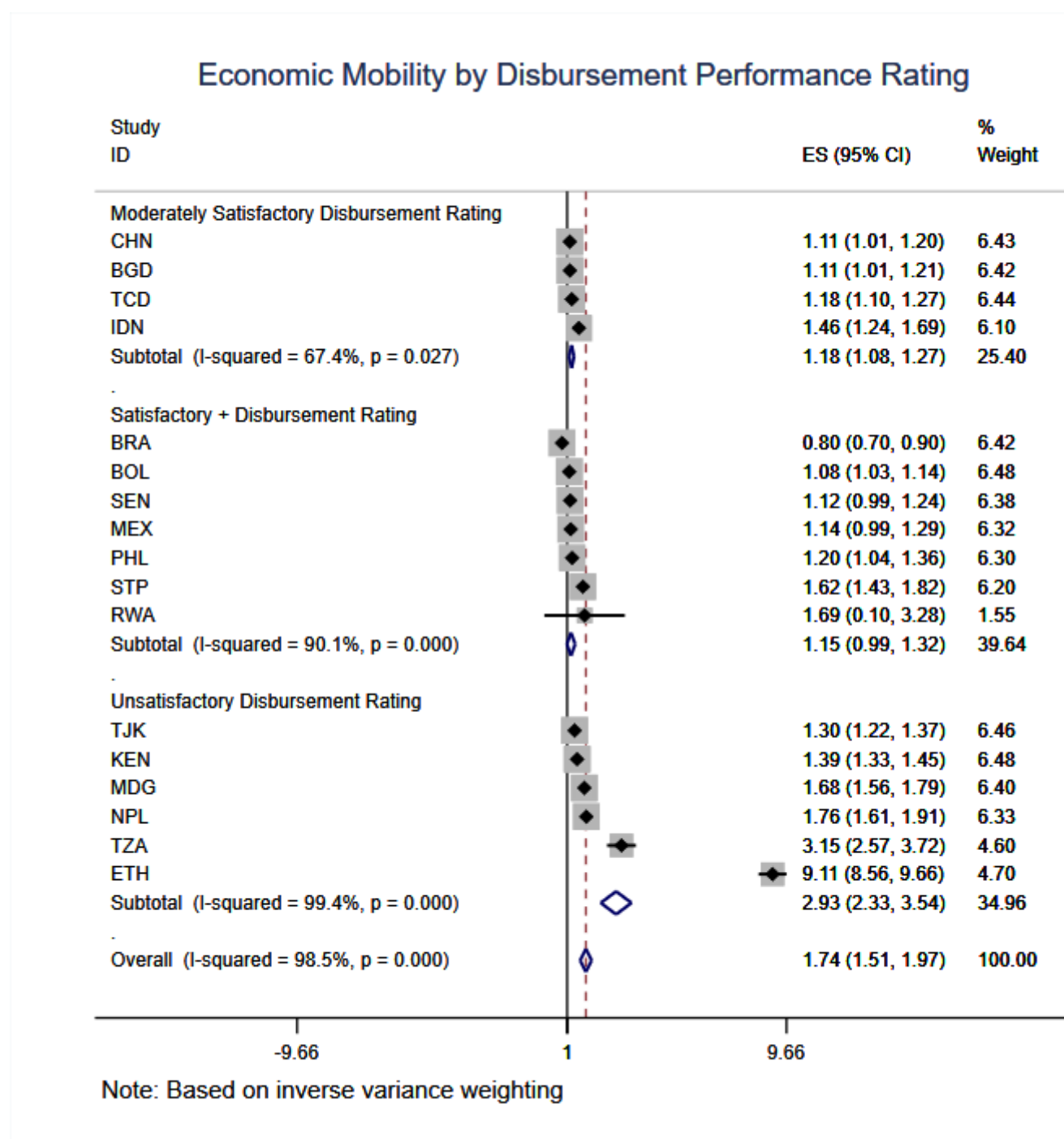
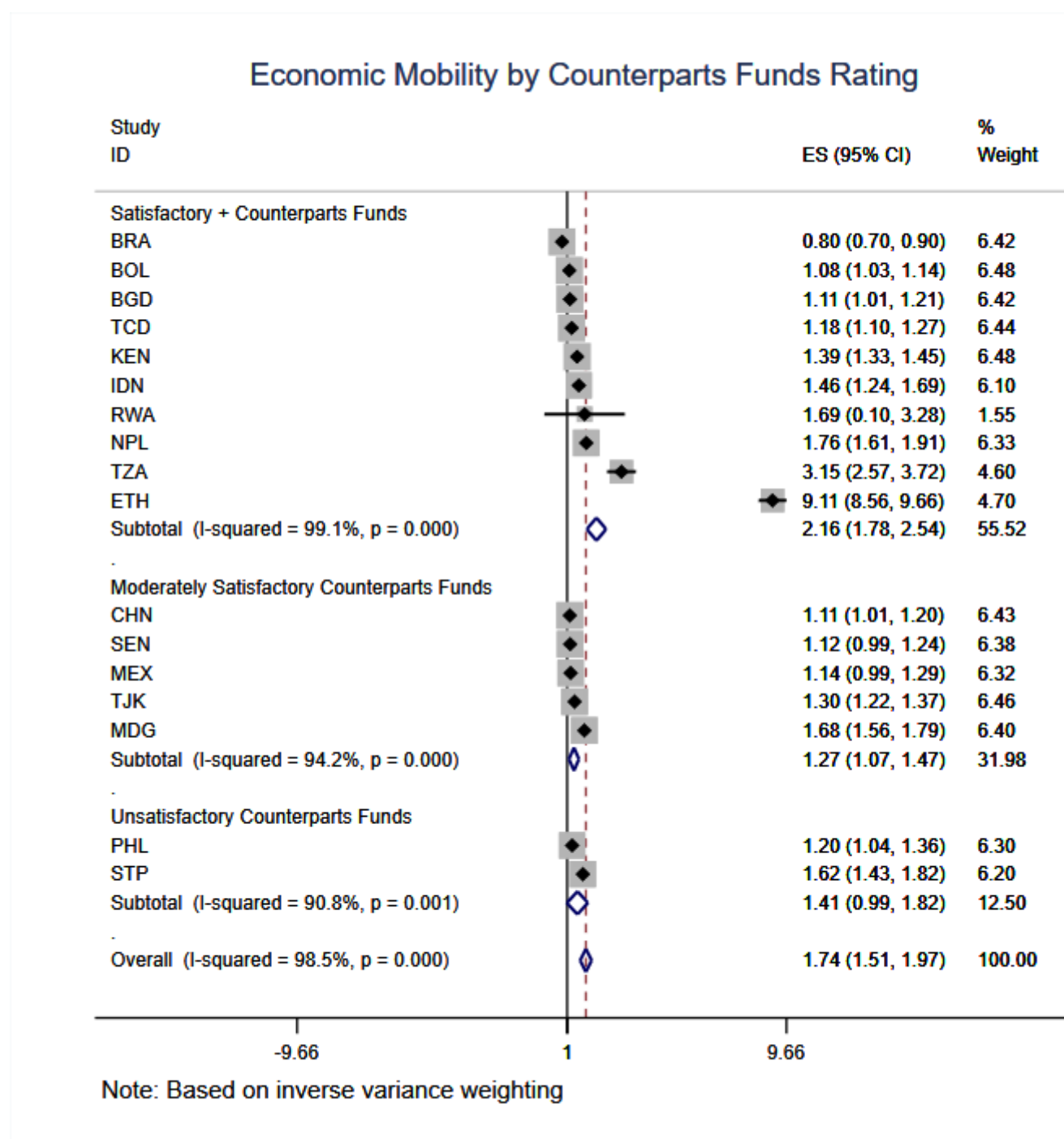
Table 31

Table 32

Appendix - Annex III

Sample correction a la Heckman

Table 33

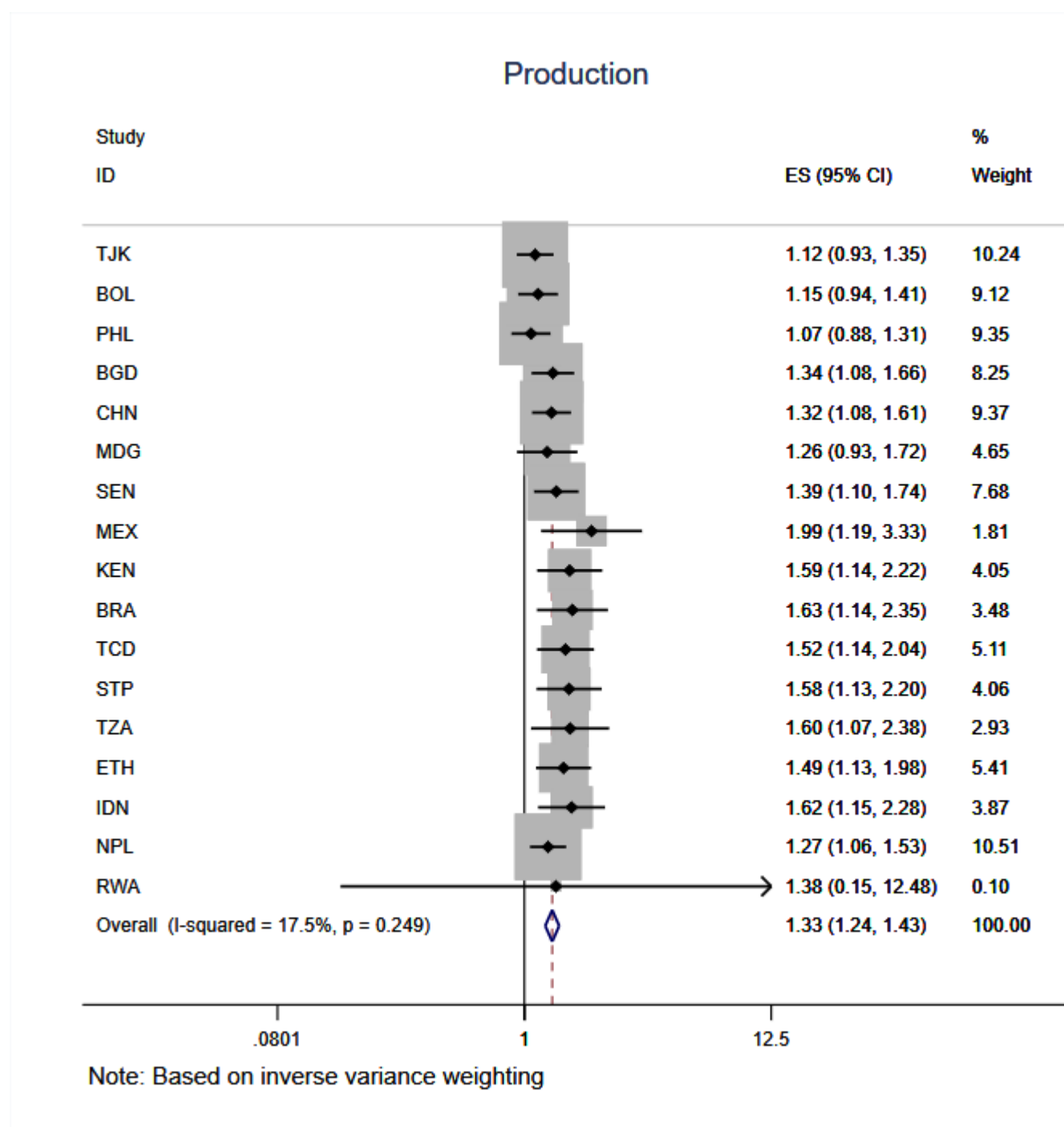


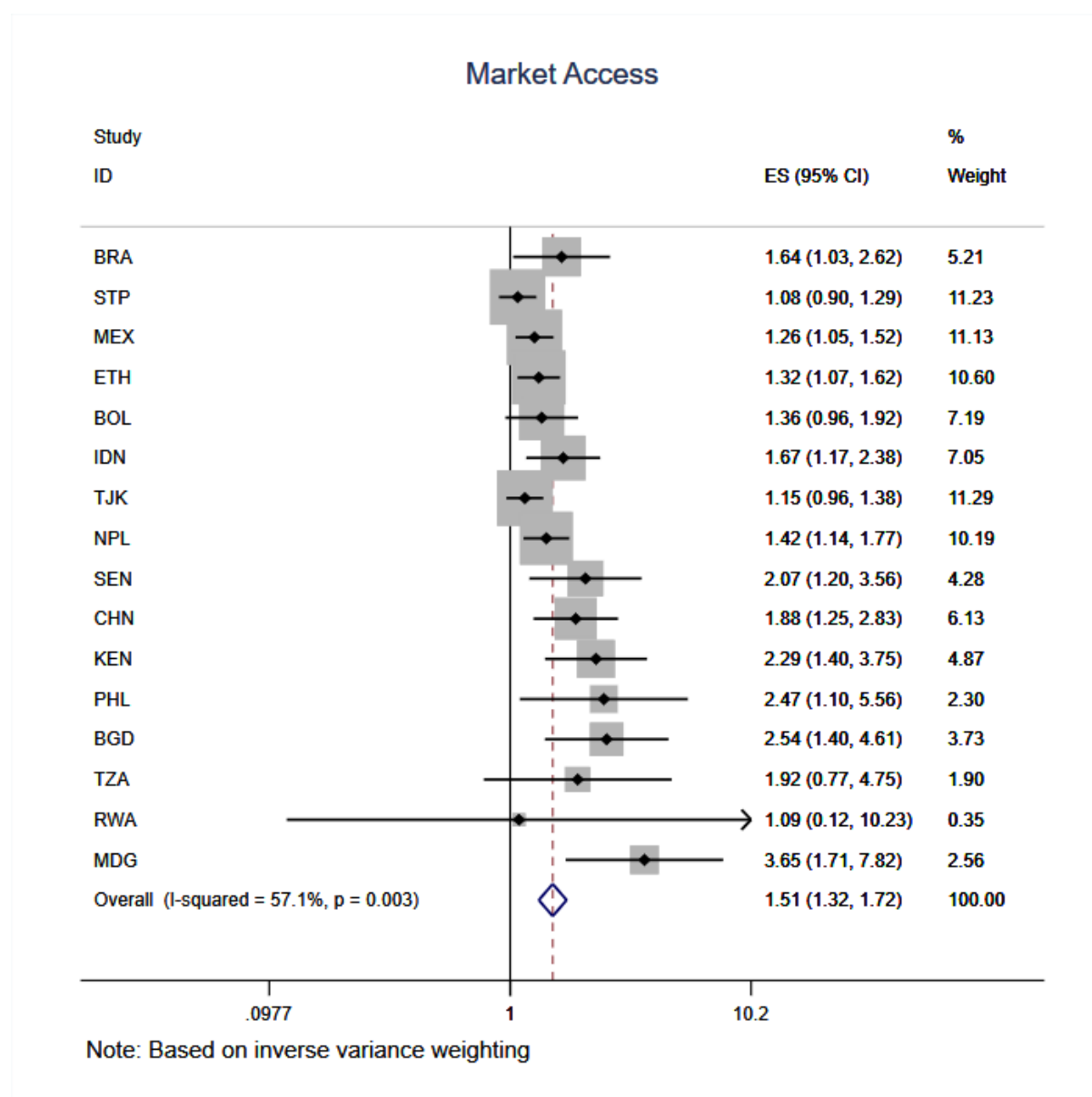
Table 34

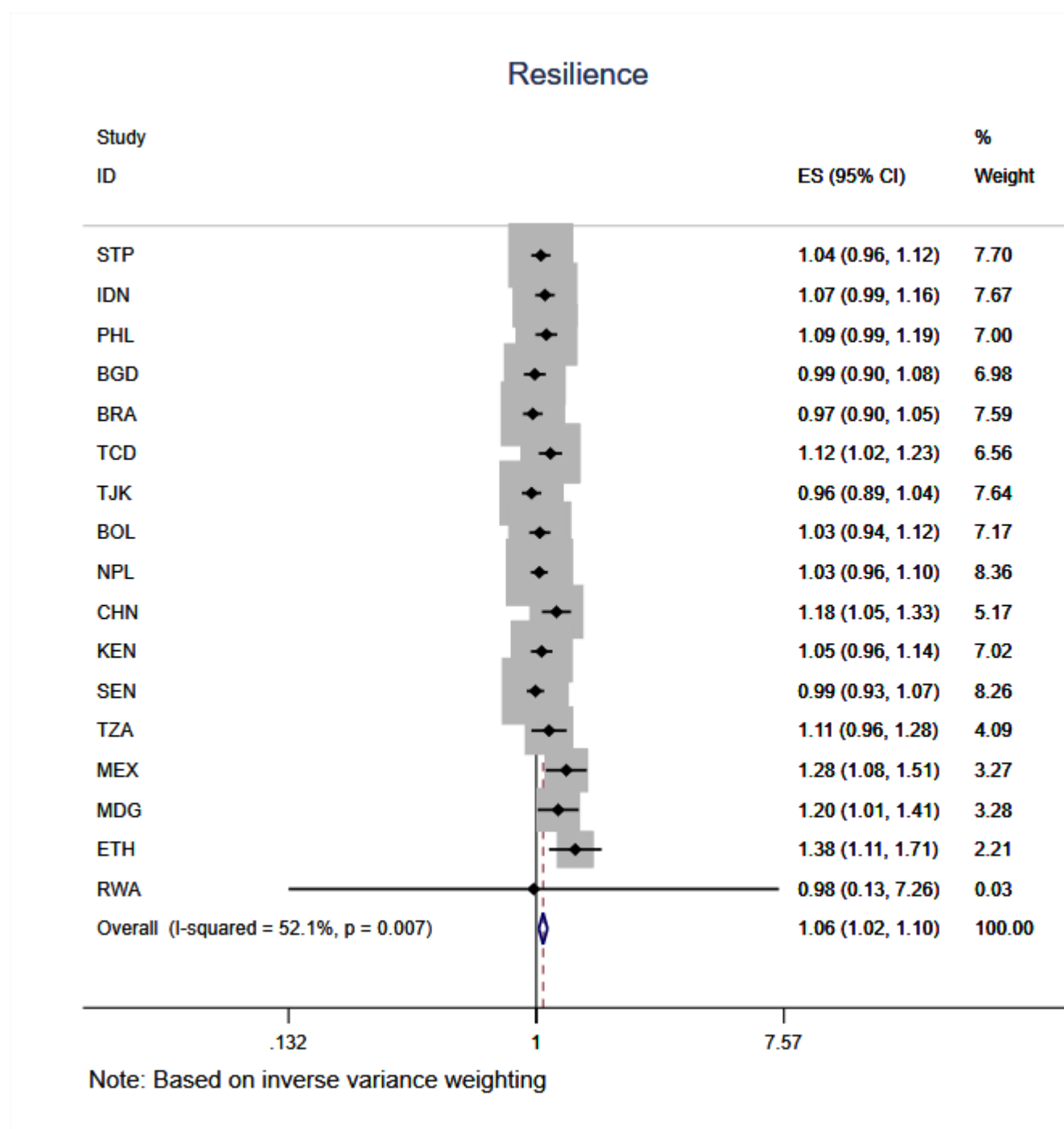
Table 35

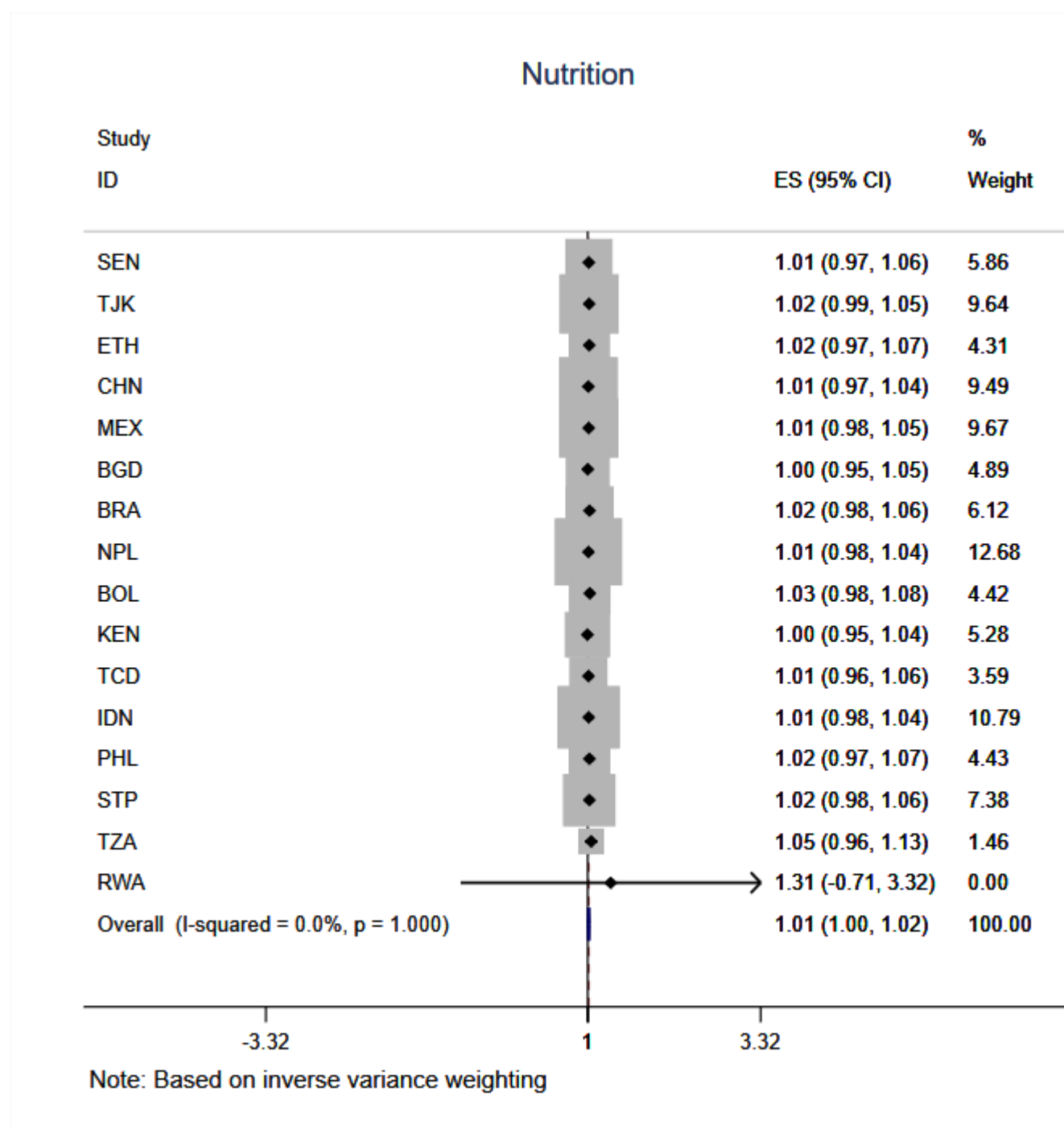
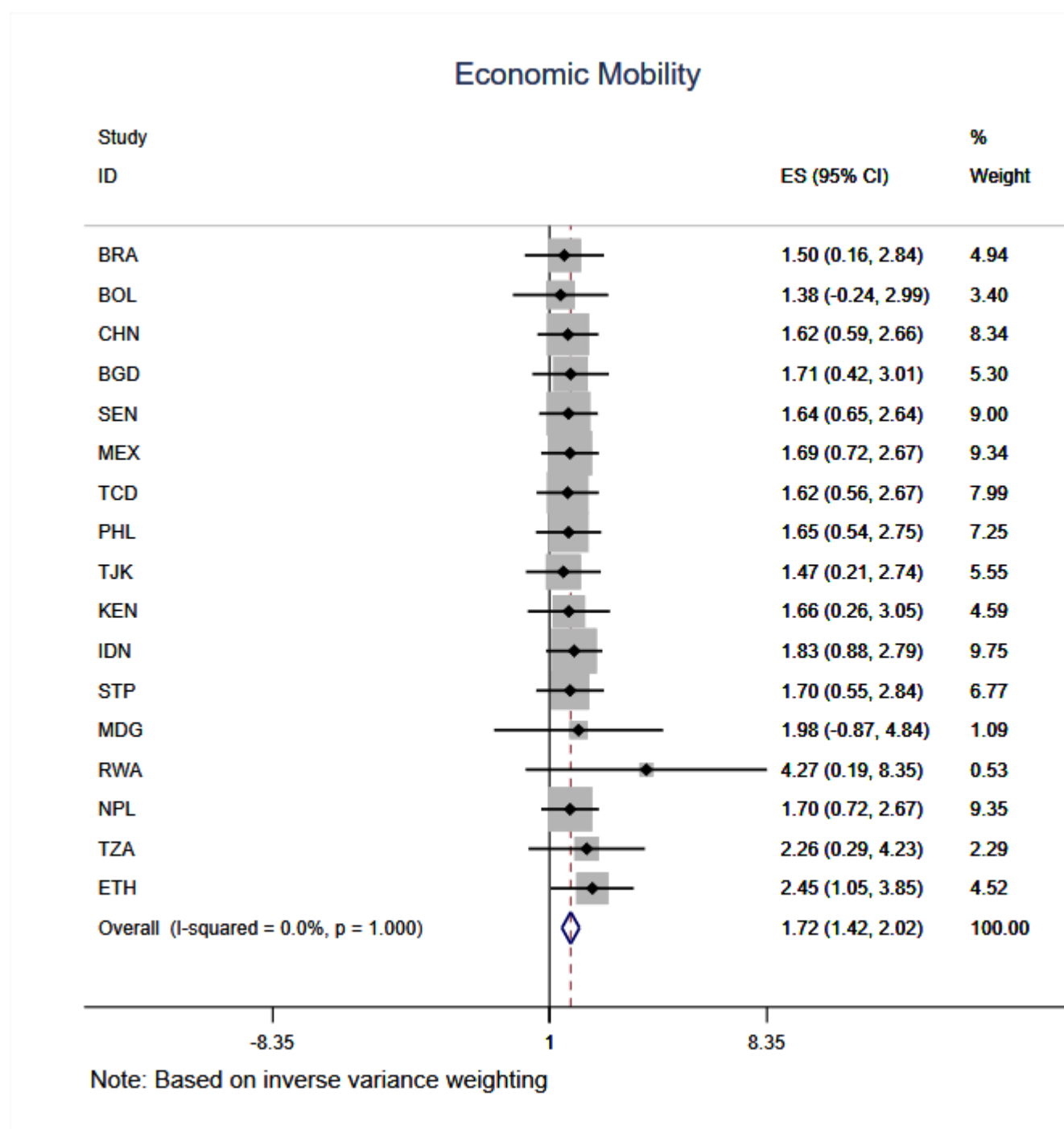
Table 36

Table 37

Appendix - Annex IV

Sensitivity Analyses results with the trim-and-fill method

The trim and fill – adjusts for bias non-parametrically. Specifically, in order to investigate for the presence of small study effects and publication bias, visual representations such as funnel or contour enhanced funnel plots are employed.

A funnel plot shows effect sizes against measures of study precision e.g. standard error. The funnel plot is employed to explore visually publication bias or more precisely small study effect. The asymmetry is evidence and maybe the result of publication bias or may be because of other reasons (heterogeneity between studies).

The contour enhanced funnel plot, can help determine whether the asymmetry of the funnel plot is due to selection bias (e.g. publication bias). The contour lines correspond to certain levels of statistical significance. Publication bias is suspect when smaller studies are absent from the non-significant regions.

Tests for funnel-plot asymmetry are useful for detecting publication bias but are not able to estimate the impact of this bias on the final meta-analysis results. The nonparametric trim-and-fill method of Duval and Tweedie (2000a, 2000b) provides a way to assess the impact of missing studies because of publication bias on the meta-analysis. It evaluates the amount of potential bias present in meta-analysis and its impact on the final conclusion. This method is typically used as a sensitivity analysis to the presence of publication bias.

Results from the Trim-and-fill method are presented in Table 38 which summarizes the original results for the meta-analysis (the observed effect size – ES) along with the imputed one from the trim and fill results (observed plus imputed ES). The full set of tables are **Error! Reference source not found.** to Table 48.

The table shows that adjusted results remain largely positive and sometimes unaltered for some domains. Bias might affect only three coefficients: Economic mobility (1.74 vs. 1.38 equivalent to 74% and 38% respectively), Market access (1.76 vs. 1.38 equivalent to 76% and 38% respectively); and Resilience indicators (1.13 vs. 1.03 equivalent to 13% and 3%, respectively). However a known limitation of Trim-and-Fill is that it can correct for publication bias that does not exist, underestimating effect sizes (Terrin et al 2003). Recommendations from recently published literature, (Simonsohn et al 2014) argued against the use of such method²³.

²³ Simonsohn, Uri and Nelson, Leif D. and Simmons, Joseph P., P-Curve and Effect Size: Correcting for Publication Bias Using Only Significant Results (April 27, 2014). Available at SSRN: <https://ssrn.com/abstract=2377290> or <http://dx.doi.org/10.2139/ssrn.2377290>

Table 38: Results from the Trim and Fill method

Production				
	N. projects	ES	Lower CI	Upper CI
Observed	17	1.44	1.26	1.65
Observed+Imputed	17	1.44	1.26	1.65
Market Access				
	N. projects	ES	Lower CI	Upper CI
Observed	16	1.76	1.45	2.14
Observed+Imputed	21	1.38	1.13	1.67
Resilience				
	N. projects	ES	Lower CI	Upper CI
Observed	17	1.13	1.02	1.25
Observed+Imputed	20	1.04	0.91	1.18
Nutrition				
	N. projects	ES	Lower CI	Upper CI
Observed	16	1.01	1	1.03
Observed+Imputed	17	1.01	1	1.03
Economic Mobility				
	N. projects	ES	Lower CI	Upper CI
Observed	17	1.74	1.51	1.97
Observed+Imputed	18	1.38	1.1	1.67

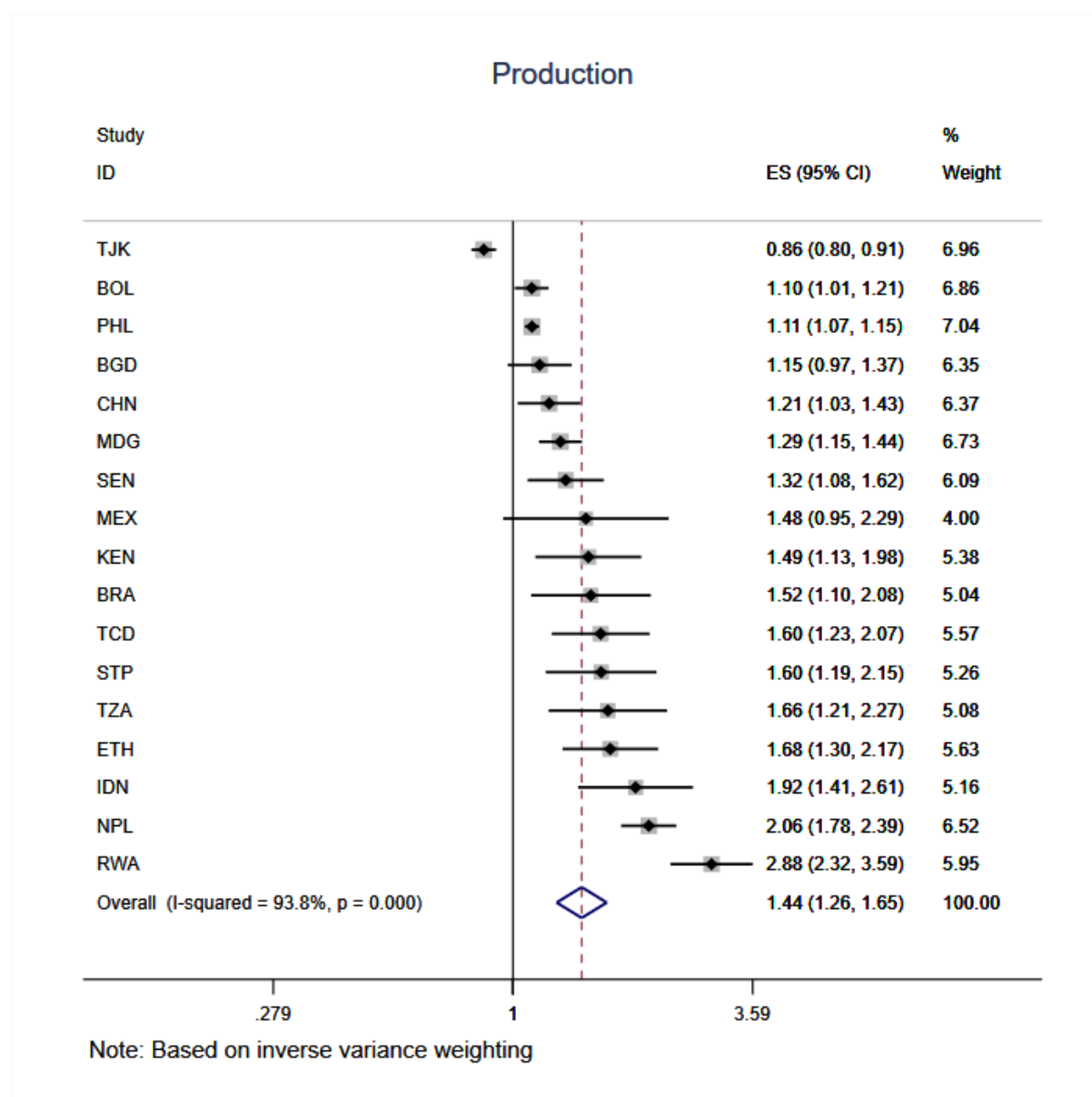
Table 39

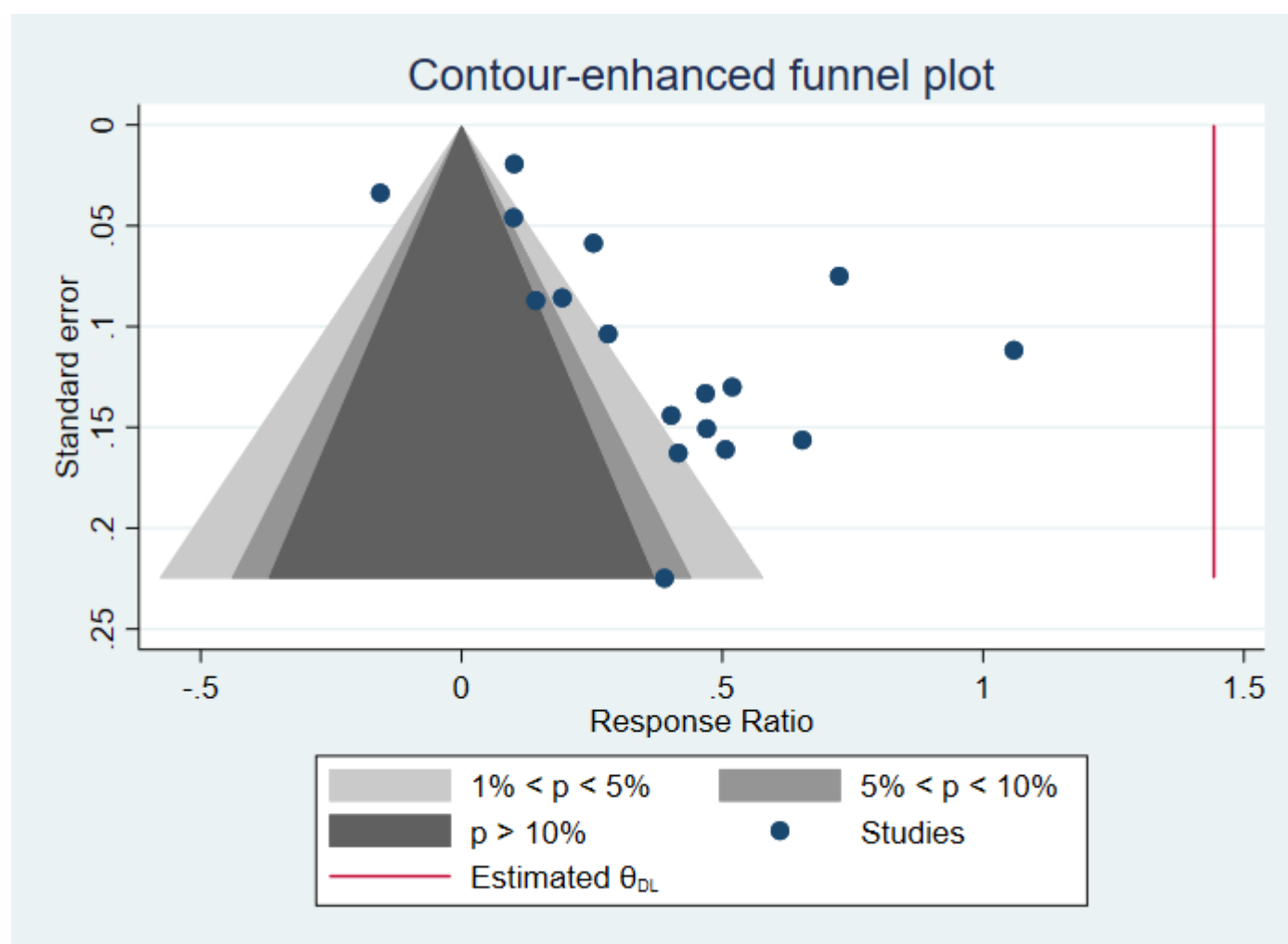
Table 40

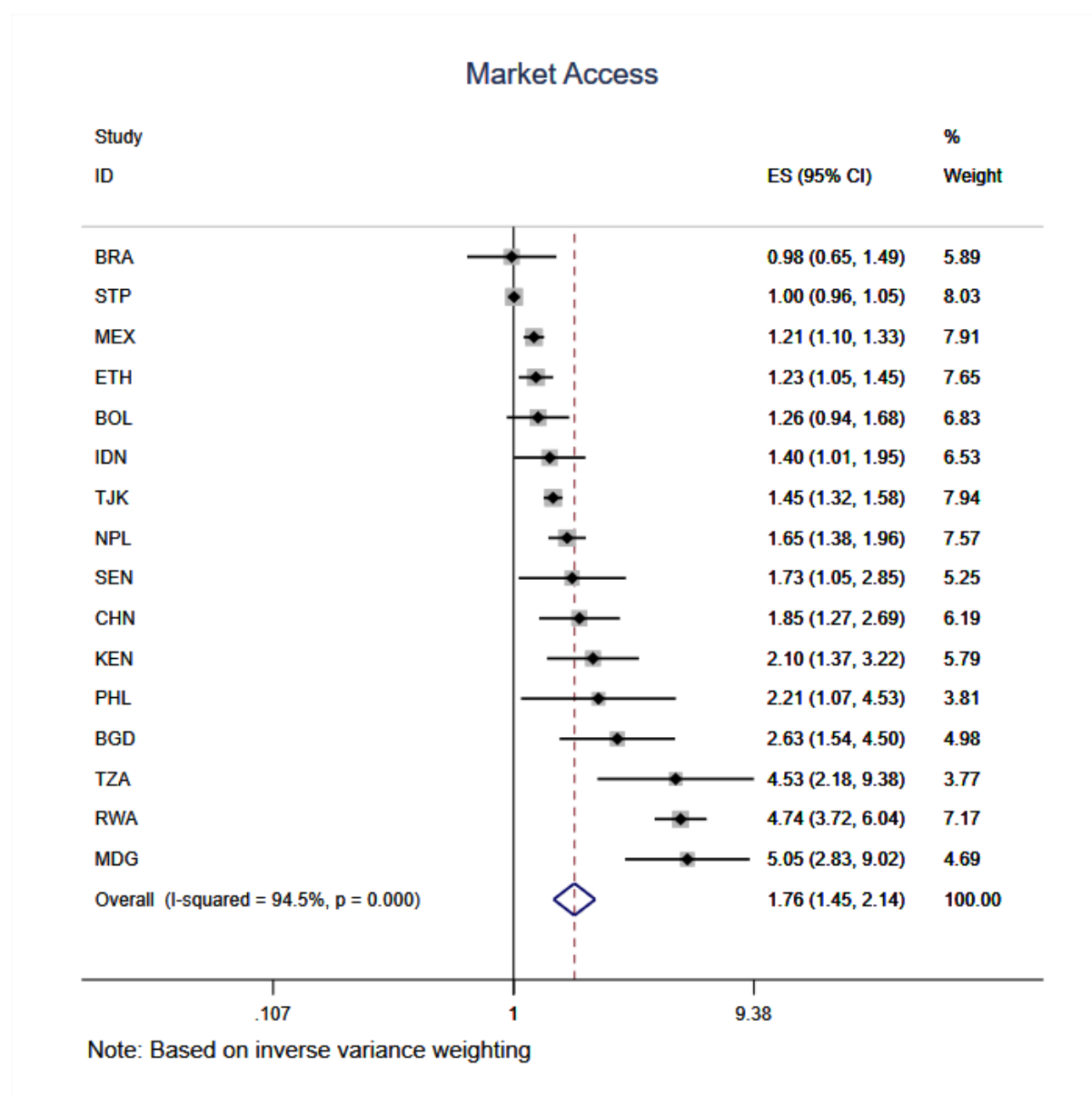
Table 41

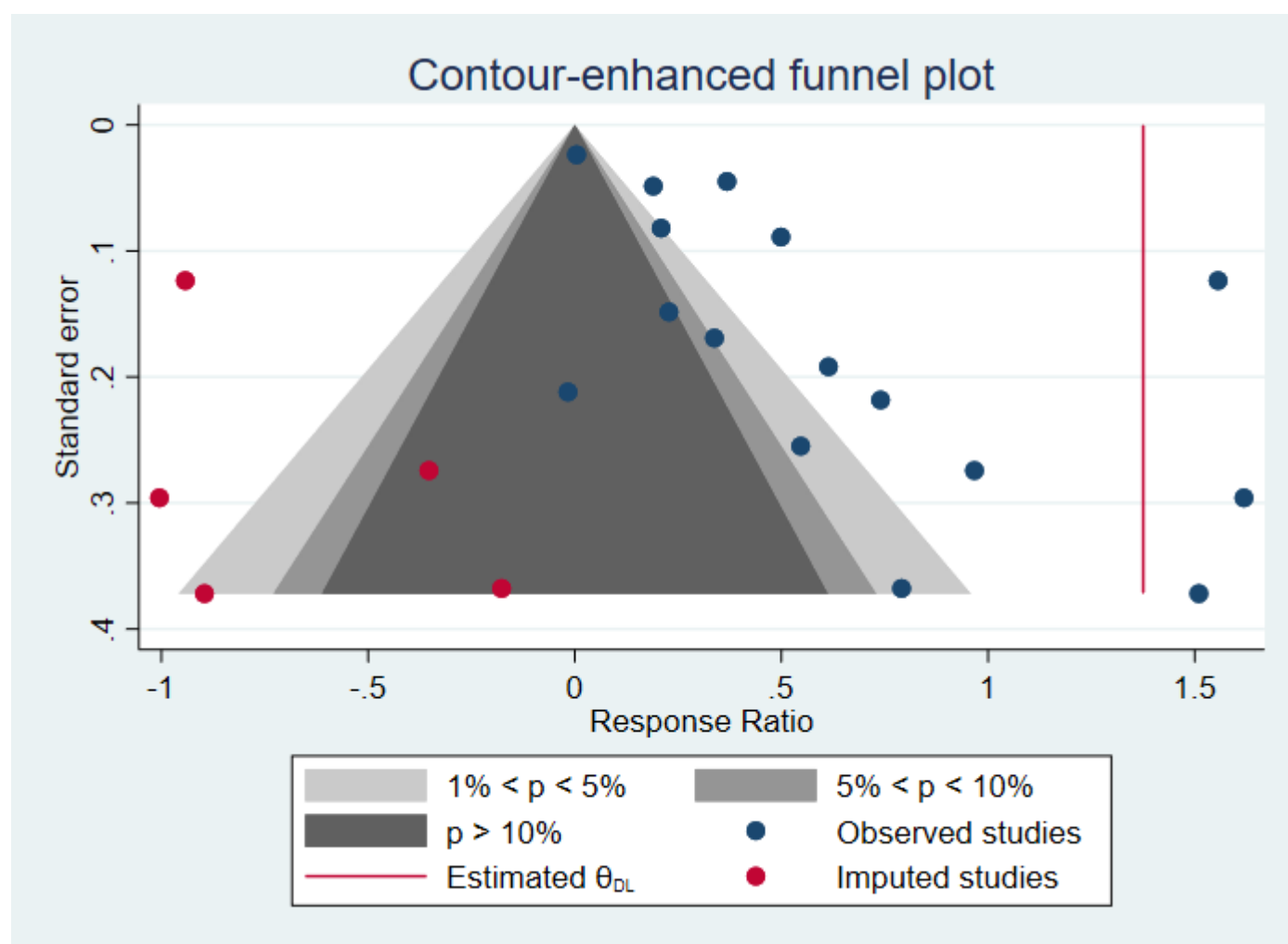
Table 42

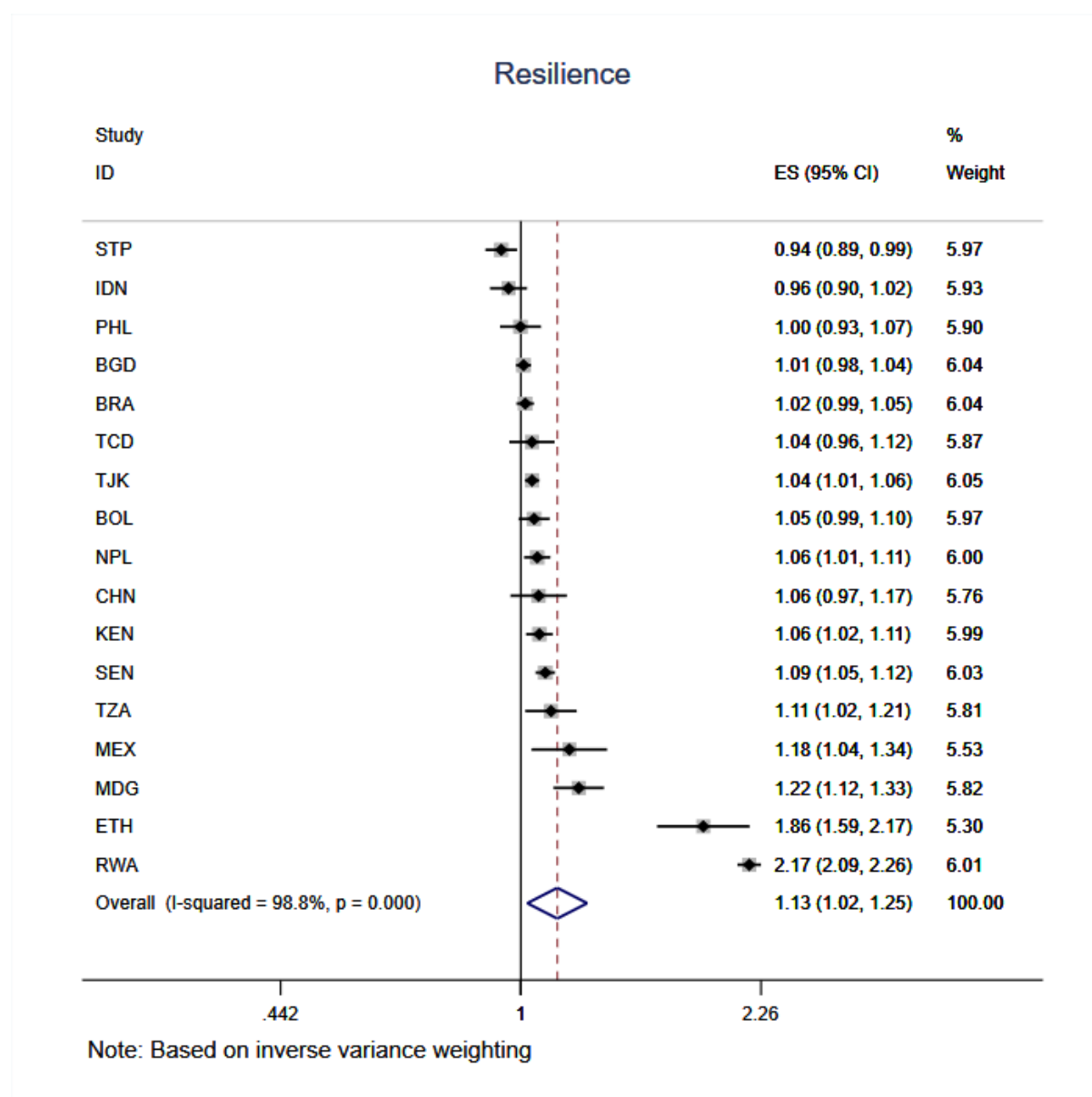
Table 43

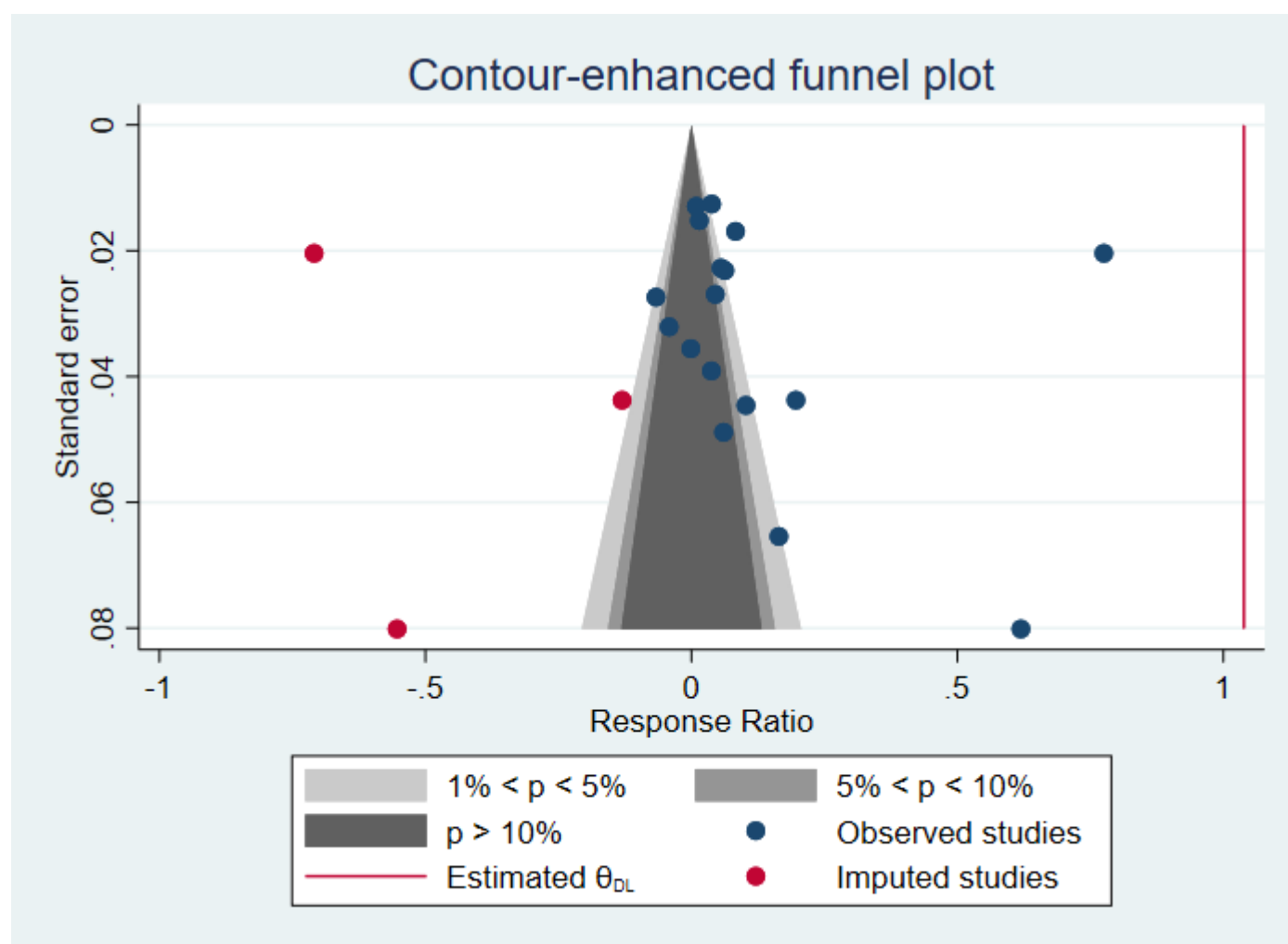
Table 44

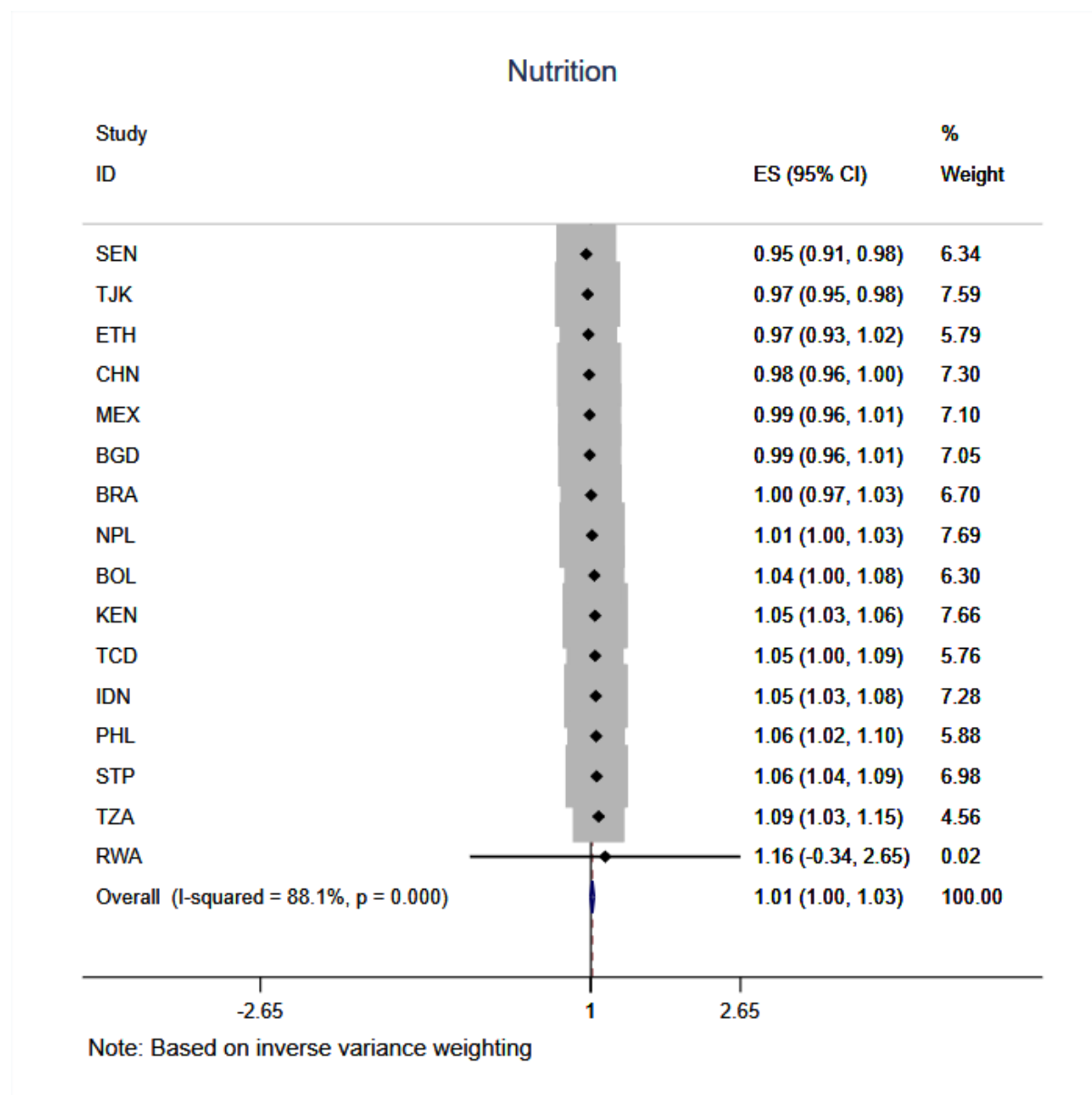
Table 45

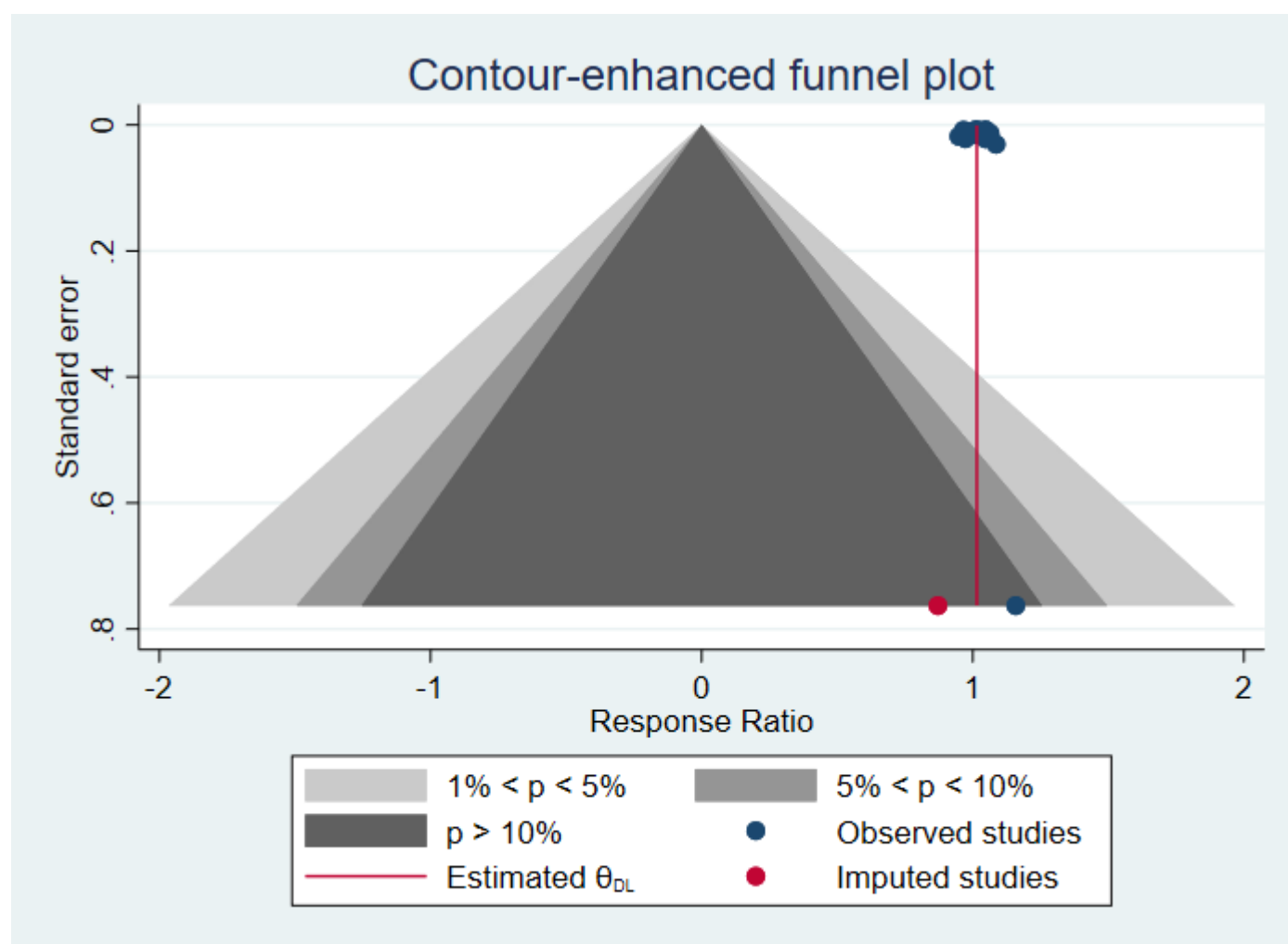
Table 46

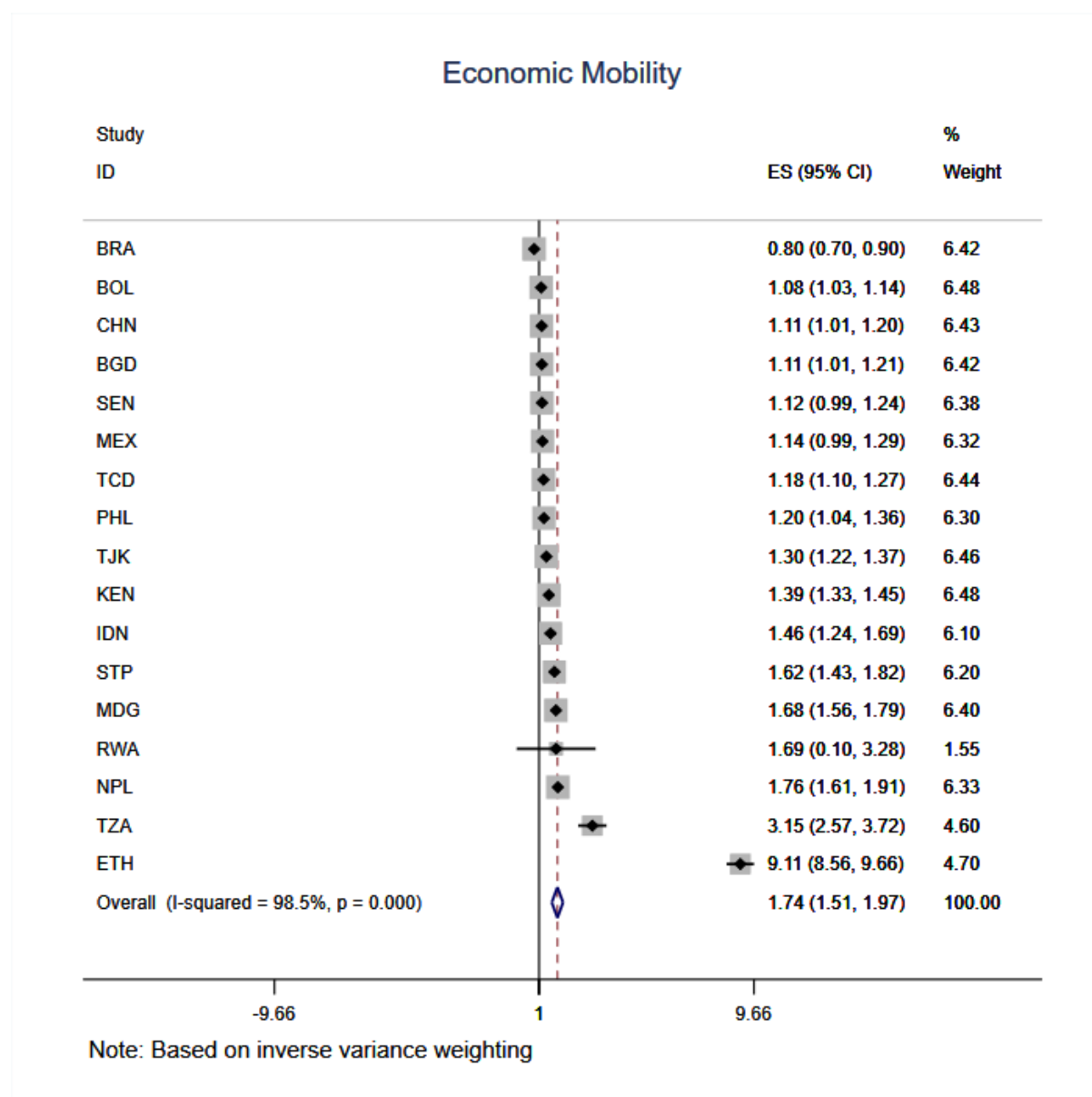
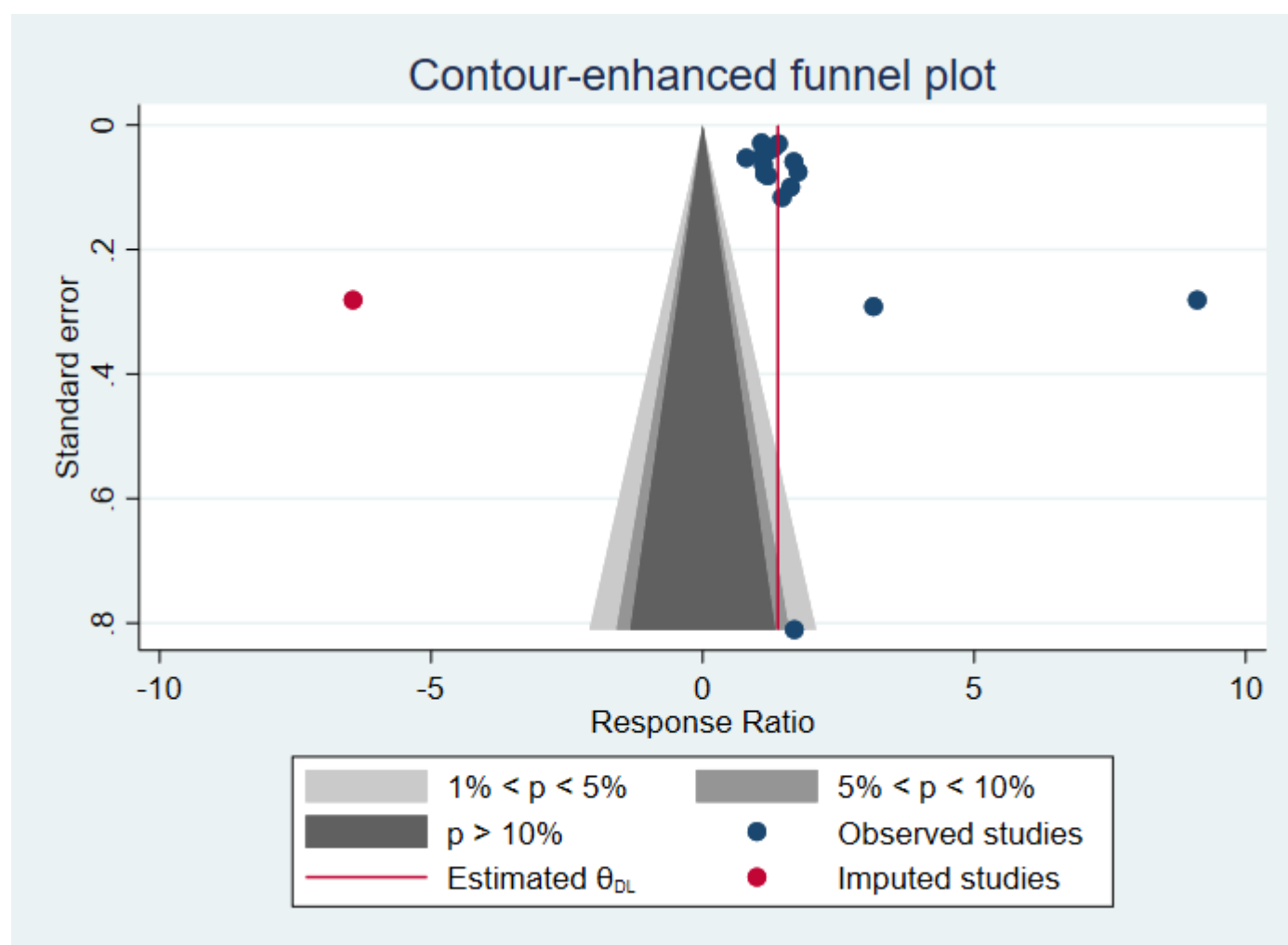
Table 47

Table 48

References

- Andrews, Isaiah and Emily Oster. 2018. "Weighting for External Validity." National Bureau of Economic Research Working Paper 23826. <http://www.nber.org/papers/w23826>.
- Banerjee, Amitav and Chaudhury, Suprakash. 2010. "Statistics without tears: Populations and samples." *Industrial Psychiatry Journal*. 19(1): 60–65.
- Bown, M.J. and A.J. Sutton. 2010. "Quality Control in Systematic Reviews and Meta-analyses." *European Journal of Cardiovascular and Endovascular Surgery*, 40: 669–677.
- Breslow NE, Day NE. 1980. *Statistical Methods in Cancer Research: Vol. I - The Analysis of Case-Control Studies*. Lyon: International Agency for Research on Cancer.
- Copas, John and Jian Qing Shi. 2000. "Meta-analysis, funnel plots, and sensitivity analysis." *Biostatistics*, 1(3): 247–262.
- Clarke, Kevin A, Randall W. Stone. 2019. "The Unobserved IMF." *The Political Economy of International Organizations*. https://www.peio.me/wp-content/uploads/2019/01/PEIO12_paper_109.pdf
- Cortes, Corinna; Mohri, Mehryar; Riley, Michael; Rostamizadeh, Afshin (2008). *Sample Selection Bias Correction Theory (PDF)*. *Algorithmic Learning Theory. Lecture Notes in Computer Science*. 5254. pp. 38–53.

- Cortes, Corinna; Mohri, Mehryar (2014). "Domain adaptation and sample bias correction theory and algorithm for regression" (PDF). *Theoretical Computer Science*. 519: 103–126.
- Delgado-Rodríguez M, Llorca J. "Bias" *Journal of Epidemiology & Community Health* 2004;58:635-641.
- Greenland, S. (2005) Multiple-bias modelling for analysis of observational data. *J. R. Statist. Soc. A*, 168, 267–291
- Heckman, James J., Sergio Urzua and Edward Vytlacil. 2006. "Understanding Instrumental Variables In Models With Essential Heterogeneity," *Review of Economics and Statistics*, 88(3, Aug): 389-432
- Henmi M, Copas JB. 2010. Confidence intervals for random effects meta-analysis and robustness to publication bias. *Statistics in Medicine* 29: 2969-2983. doi: 10.1002/sim.4029
- Lin L (2018) Bias caused by sampling error in meta-analysis with small sample sizes. *PLoS ONE* 13(9): e0204056. <https://doi.org/10.1371/journal.pone.0204056>
- Mathur M.B & VanderWeele Tyler J. (2019): Sensitivity Analysis for Unmeasured Confounding in Meta-Analyses, *Journal of the American Statistical Association*, DOI: 10.1080/01621459.2018.1529598
- Mavridis D, Salanti G. How to assess publication bias: funnel plot, trim-and-fill method and selection models. *Evidence-Based Mental Health* 2014;17:30.
- Oster, Emily. 2019. "Unobservable Selection and Coefficient Stability: Theory and Evidence." *Journal of Business & Economic Statistics*, 37(2):187-204.
- Shi, Linyu and Lifeng Lin. 2019. "The trim-and-fill method for publication bias: practical guidelines and recommendations based on a large database of meta-analyses." *Medicine*, 98(23): 1-11.
- Stuart, Elizabeth A., Stephen R. Cole, Catherine P. Bradshaw, and Philip J. Leaf. 2011. "The use of propensity scores to assess the generalizability of results from randomized trials." *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 174(2): 369-386.