

Online appendix C: Risk of bias assessment tool

The following table provides a provisional tool to guide the risk of bias assessment for quantitative impact evaluations. If necessary, we could amend the tool to better inform the appraisal of primary studies.

Provisional risk of bias assessment tool (RCT)

General	ID	EPPI ID		
General	Study first author	Open answer		
General	Time taken to complete assessment	Minutes		
General	Design type: What type of study design is used?	1= Randomised controlled trial (RCT) (random assignment to households/individuals) or quasi-RCT 2= Cluster-RCT (quasi-RCT)	-	
General	Methods used for analysis: Which methods are used to control for selection bias and confounding?	1 = Statistical matching (PSM, CEM, covariate matching) 2 = Difference in differences (DID) estimation methods 3 = IV-regression (2-stage least squares or bivariate probit) 4 = Heckman selection model 5 = Fixed effects regression 6 = Covariate adjusted estimation 7 = Propensity weighted regression 8 = Comparison of means 9 = Other (please state)	-	
General	Design and analysis method description	Open answer	Briefly describe the study design and analysis method undertaken by the authors.	

General	Study population	Open answer	Provide any details in the paper that describe how the study population was selected, covering: a) How is the population selected? what is the sampling strategy to recruit participants from that population into the study? b) What are the characteristics of that study participants? c) Was this a pilot program aimed at being scaled up? d) Were there specific factors of success or failure in the implementation?	
General	Type of comparison group	1=No intervention (service delivery as usual) 2=Other intervention 3=Pipeline (wait-list) control (still service delivery as usual)	Indicate type of comparison group	
General	Type of comparison group (if other)	Open answer		
General	Ethical clearance	Open answer	Provide any details of ethical research clearances granted. Report unclear if this information is not available.	
General	Study registration	Open answer	Provide any details of study registration, including registry IDs, etc.	
1: Assignment mechanism - Assessment	Assignment mechanism: Was the allocation or identification mechanism random or as good as random?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	a) The authors describe a random component in sequence generation/ randomization method (e.g. lottery, coin toss, random number generator) and assignment is performed for all units at the start of the study centrally or using a method concealed from participants and intervention delivery b) If public lottery is used for the sequence	Score "Yes" if all criterion a), b), c) and d) are satisfied. Score "Probably Yes" if only criterion a) and b) are not satisfied OR if only criteria c) is not satisfied. Score "Unclear" if d) is not satisfied because no balance table is

			<p>generation, authors provide detail on the exact settings and participants attending the lottery.</p> <p>c) If a special randomization procedure is used to ensure balance, it is well described and justified given the study setting (stratification, pairwise matching, unique random draw, multiple random draws etc).</p> <p>d) A balance table is reported suggesting that allocation was random between all groups including subgroup receiving different treatment within control or treatment groups (if the comparison is relevant for this assessment).</p>	<p>reported.</p> <p>Score "Probably No" if d) is not satisfied because there is no balance table reported and there is evidence suggesting a problem in the randomization, such as baseline coefficients in a diff-in-diff regression table are very different or sample size is too small for the procedure used (using stratification when there are less than two units for each intervention and control group in each strata can lead to imbalance).</p> <p>Score "No" if d) is not satisfied because there are large imbalances concerning a large number of variables, providing evidence that the assignment was not random. If this is scored as no, use the NRS tool.</p>
1: Assignment mechanism - Justification	Assignment justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	
2: Unit of analysis - Assessment	Unit of analysis: Is unit of analysis in cluster allocation addressed in standard error calculation ?	1=Yes 2=No 3=Not reported/unclear 4=Not applicable	<p>Score "Yes" if UoA = UoR OR if UoA != UoR and standard errors are clustered at the UoR level OR data is collapsed to the UoR level</p> <p>Score "Not reported/unclear" if not enough information is provided on the way the standard errors were calculated or what the</p>	

			<p>unit of analysis is.</p> <p>Score "Not applicable" if it is not a cluster RCT.</p> <p>Score "No" otherwise.</p>	
2: Unit of analysis - Justification	Method used to address differences between UoA and unit of data collection	Open answer		
3: Selection bias - Assessment	Selection bias Was any differential selection into or out of the study (attrition bias) adequately resolved?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	<p>Score "Yes" if there is no attrition or attrition falls into the green zone and the study establishes that attrition is randomly distributed (e.g. by presenting balance by key characteristics across groups) AND if survey respondents were randomly sampled.</p> <p>Score "Probably yes" if attrition falls into the green zone AND if survey respondents were randomly sampled.</p> <p>Score "Unclear" if there is an attrition problem but no information provided on the relationship between attrition and treatment status, OR if there is not enough information on how the population surveyed was sampled.</p> <p>Score "Probably no" if there is attrition which is likely to be related to the intervention OR there is some indication that the survey respondents were purposely</p>	

			<p>sampled in a way that might have led the sampling to be different between treatment and control groups, or attrition falls into the yellow zone.</p> <p>Score "No" if attrition falls into the red zone.</p>	
3: Selection bias - Justification	Selection bias justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	
3: Confounding - Assessment	Confounding and group equivalence: Was the method of analysis executed adequately to ensure comparability of groups throughout the study and prevent confounding	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	<p>a) Baseline characteristics are similar in magnitude;</p> <p>b) Unbalanced covariates at the individual and cluster level are controlled in adjusted analysis;</p> <p>c) Adjustments to the randomization were taken into account in the analysis (stratum fixed effects, pairwise matching variables)? (Bruhn and McKenzie 2009)</p>	<p>Score "Yes" if criterion a) and b) are satisfied;</p> <p>Score "Probably yes" if a) is not satisfied but b) is satisfied and imbalances are small in magnitude OR if only a) is satisfied.</p> <p>Score "Unclear" if no balance table is provided or if imbalances are controlled for but they are very large in magnitude and assignment mechanism is not coded as "Yes" or "Probably yes"</p> <p>Score "Probably no" if a) and b) are not satisfied and the magnitude of imbalances are small</p> <p>Score "No" if a) and b) are not satisfied and the magnitude of imbalances are large and covariates are clear determinant of the outcomes.</p>

3: Confounding - Justification	Confounding justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	
--------------------------------------	---------------------------	-------------	---	--

<p>4: Deviations from intended interventions - Assessment</p>	<p>Deviations from intended interventions: Spill-overs, cross-overs and contamination: was the study adequately protected against spill-overs, cross-overs and contamination?</p>	<p>1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear</p>	<p>a) There was no implementation issues that might have led the control participants to receive the treatment (implementer's mistake). b) The intervention is unlikely to spill-over to comparisons (e.g. participants and non-participants are geographically and/or socially separated from one another and general equilibrium effects are not likely) or the potential effects of spill overs were measured (e.g. variation in the % of unit within a cluster receiving the treatment). c) There is no risk of contamination by external programs: the treatment and comparisons are isolated from other interventions which might explain changes in outcomes. d) There is nothing in the surveys that might have given the control participants an idea of what the other group might receive OR they did but there is no risk that this has changed their behaviours; AND the survey process did not reveal information to the control group that they did not have before (e.g. the study aims to measure increase in take up of a service or product that participants might not know about) Authors might put something in place in the design of the study that allows to control for that survey effect (e.g. a pure control with no monitoring except baseline end line)</p>	<p>Score "Yes" if criterion a), b), c) and d) are satisfied;</p> <p>Score "Probably yes" if there is no obvious problem but there is no information reported on potential risks related to spill overs, contamination, or survey effects in the control group OR if there were issues with spill-overs but they were controlled for or measured.</p> <p>Score "Unclear" if spill-overs, cross-overs, survey effects and/or contamination are not addressed clearly.</p> <p>Score "Probably no" if any of the criterion a), b), c) or d) are not satisfied but the scale of the issue is not clear.</p> <p>Score "No" if any of the criterion a), b), c) or d) are not satisfied and happened at a large scale in the study.</p>
---	---	---	--	---

4: Deviations from intended interventions - Justification	Deviations justification	Open answer	<p>Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).</p> <p>For example, intervention groups are geographically separated, authors use intention to treat estimation or instrumental variables to account for non-adherence, and survey questions are not likely to expose individuals in the control group to information about desirable behaviours ('survey effects').</p>	
5. Performance bias - Assessment	Performance bias: Was the process of monitoring individuals unlikely to introduce motivation bias among participants?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	<p>a) The authors state explicitly that the process of monitoring the intervention and outcome measurement is blinded and conducted in the same frequency for treatment and control groups, or argue convincingly why it is not likely that being monitored could affect the performance of participants in treatment and comparison groups in different ways (such as resulting in Hawthorne or John Henry effects).</p> <p>b) The outcome is based on data collected in the context of a survey, and not associated with a particular intervention trial, or data are collected from administrative records or in the context of a retrospective (ex post) evaluation.</p>	<p>Score "Yes" if either criterion a) or b) are satisfied;</p> <p>Score "Probably yes" if the study is based on data collected during a trial and there is no obvious issue with the monitoring processes but authors do not mention potential risks.</p> <p>Score "Unclear" if it is not clear whether the authors use an appropriate method to prevent Hawthorne and John Henry Effects (e.g. blinding of outcomes and, or enumerators, other methods to ensure consistent monitoring across groups). Hawthorne effects may result where participants know that they are being observed and John Henry Effects may result from participant knowledge of being</p>

				<p>compared.</p> <p>Score "Probably no" if there was imbalance in the frequency of monitoring in intervention groups, which might have influenced participants' behaviours.</p> <p>Score "No" if neither criterion a) or b) are satisfied.</p>
5. Performance bias - Justification	Performance bias justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	

6. Outcome measurement bias - Assessment	Outcome measurement bias: Was the study free from biases in outcome measurement?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	<p>a) Outcome assessors are blinded or the outcome measures are not likely to be biased by their judgement.</p> <p>b) For self-reported outcomes: respondents in the intervention group are not more likely to have accurate answers due to recall bias;</p> <p>c) For self-reported outcomes: respondents do not have incentives to over/under report something related to their performance or actions, OR researchers put in place mechanisms to reduce the risk of reporting bias (researchers not strongly involved in the implementation of the program and it is clear that their answers to the survey will not affect what they receive in the future) OR authors have measured the risks of bias through falsification tests or measuring the effect on placebo outcomes in cases where there was a risk of reporting bias.</p> <p>d) Timing issue: the data collection period did not differ between intervention and comparison group, the baseline data is not likely to be affected by the beginning of the intervention or affects a small percentage of the study participants.</p>	<p>Score "Yes" if criterion a), b), c) and d) are satisfied:</p> <p>Score "Probably yes" if there is a small risk related to any of a), b), c) or d) and there is no more information provided to justify the absence of bias OR if there was a high risk of bias but authors have either controlled it in their design or measured it with a placebo outcomes.</p> <p>Score "Unclear" if there is a high risk related to any of a), b), c) or d) and there is no more information provided to justify the absence of bias.</p> <p>Score "Probably no" if there are high risk related to a), b), c) or d) and it is clear that authors were not able to control for this bias.</p> <p>Score "No" if there is evidence of bias.</p>
6. Outcome measurement bias - Justification	Outcome measurement justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	

7. Reporting bias - Assessment	Analysis reporting: Was the study free from selective analysis reporting?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	<p>a) A pre-analysis plan or trial protocol is published and referred to or the trial was pre-registered or the outcomes were pre-registered;</p> <p>b) Authors report results corresponding to the outcomes announced in the method section (there is no outcome reporting bias);</p> <p>c) Authors report results of unadjusted analysis and intention to treat (ITT) estimation, alongside any adjusted and treatment-on-the-treated/complier-average-causal-effects analysis.)</p> <p>d) Authors use the appropriate analysis method (use baseline data when available) and different treatment arms are differentiated in the analysis</p> <p>e) Authors have reported all the analysis which could help understand the results and no other bias is assessed as unclear due to the lack of an important analysis (e.g. a balance table or a subgroup analysis)</p>	<p>Score "Yes" if all the criterion a), b), c), d), and e) are satisfied;</p> <p>Score "Probably yes" if all the conditions are met except a), or if all the conditions are met but there is some element missing that could have helped understand the results better (e);</p> <p>Score "Unclear" if there is not enough information to determine that there is an analysis missing;</p> <p>Score "Probably no" if any of the criterion b), c) or d) are not satisfied;</p> <p>Score "No" if any of the criterion b), c) or d) are not satisfied and there is evidence that the analysis results would be different because large imbalances were not controlled for, compliance was very low and ITT estimation was not reported or different treatment arms were pooled.</p>
8. Reporting bias - Justification	Analysis reporting justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	
9. Other bias - Assessment	Other risks of bias Is the study free from other sources of bias?	1= Yes, 4 = No		
9. Other bias - Justification	Other bias justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages). For	

			example, information is collected using a different survey instrument in different intervention groups; measurement of the intervention received is unclear.	
10. Blinding - observers - Assessment	Blinding of participants?	1=Yes 2=No 8=unclear 9= N/A	If there is no information, code NO. If there is information but it is ambiguous, code UNCLEAR.	
10. Blinding - observers - Assessment	Blinding of outcome assessors?	1=Yes 2=No 8=unclear 9= N/A	If there is no information, code NO. If there is information but it is ambiguous, code UNCLEAR.	
10. Blinding - analysts - Assessment	Blinding of data analysts?	1=Yes 2=No 8=unclear 9= N/A	If there is no information, code NO. If there is information but it is ambiguous, code UNCLEAR.	
10. Blinding - method(s)	Method(s) used to blind	Open answer (including describe method of placebo control)No 9= N/A	Describe method(s) used to blind	
11. External validity - Assessment	External validity	Open answer	a) What do authors say about external validity?	Include all information that can help assess the external validity of the results.

Provisional risk of bias assessment tool (QED)

General	ID	EPPI ID		
General	Time taken to complete assessment	Minutes		
General	Study first author	Open answer		
General	Outcome	Open answer		

General	Study design: What type of study design is used?	<p>1= Natural experiment: randomised or as-if randomised</p> <p>2= Natural experiment: regression discontinuity (RD)</p> <p>3= CBA (non-randomised assignment with treatment and contemporaneous comparison group, baseline and end line data collection) – individual repeated measurement</p> <p>4= CBA pseudo panel (repeated measurement for groups but different individuals)</p> <p>5= Interrupted time series (with or without contemporaneous control group)</p> <p>6= Panel data, but no baseline (pre-test)</p> <p>7 = Comparison group with end line data only</p>		
General	Methods used for analysis: Which methods are used to control for selection bias and confounding?	<p>1 = Statistical matching (PSM, CEM, covariate matching)</p> <p>2 = Difference in differences (DID) estimation methods</p> <p>3 = IV-regression (2-stage least squares or bivariate probit)</p> <p>4 = Heckman selection model</p> <p>5 = Fixed effects regression</p> <p>6 = Covariate adjusted estimation</p> <p>7 = Propensity weighted regression</p> <p>8 = Comparison of means</p>	-	

		9 = Other (please state)		
General	Study population	Open answer	Provide any details in the paper that describe how the study population was selected, covering: a) How is the population selected? what is the sampling strategy to recruit participants from that population into the study? b) What are the characteristics of that study participants? c) Was this a pilot program aimed at being scaled up? d) Were there specific factors of success or failure in the implementation?	
General	Ethical clearance	Open answer	Provide any details of ethical research clearances granted. Report unclear if this information is not available.	
General	Study registration	Open answer	Provide any details of study registration, including registry IDs, etc.	

1: Selection bias - Assessment	1 - Mechanism of assignment: was the allocation or identification mechanism able to control for selection bias?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear		
1: Selection bias - Justification	For regression discontinuity designs	Open answer	<p>a) Allocation is made based on a pre-determined discontinuity on a continuous variable (regression discontinuity design) and blinded to participants or;</p> <p>b) if not blinded, individuals reasonably cannot affect the assignment variable in response to knowledge of the participation decision rule;</p> <p>c) and the sample size immediately at both sides of the cut-off point is sufficiently large to equate groups on average.</p>	<p>Score "Yes" if criteria a), b), c) are all satisfied</p> <p>Score "Probably Yes" if there are minor differences in between both sides of the cut-off point but authors convincingly argue that the differences are unlikely to affect the outcome, OR individuals are not blinded and there are low risk of them affecting the assignment but the authors do not mention it.</p> <p>Score "Unclear" if it is unclear whether participants can affect it in response to knowledge of the allocation mechanism.</p> <p>Score "Probably No" if there are differences between individuals on both sides of the cut-off point, and there are doubts that the differences are due to individuals altering the assignment OR the participants are blinded but there is evidence that the decisions that determined the discontinuity is based on differences between the two groups or differences in time.</p> <p>Score "No" if the sample size is not</p>

				sufficient OR there is evidence that participants altered the assignment variable prior to assignment. If the research has serious concerns with the validity of the assignment process or the group equivalence completely fails, we recommend assessing risk of bias of the study using the relevant questions for the appropriate methods of analysis (cross-sectional regressions, difference-in-difference, etc.) rather than the RDDs questions.
1: Selection bias - Justification	For assignment based non-randomised programme placement and self-selection (studies using a matching strategy or regression analysis, excluding IV)	Open answer	<p>a) Participants and non-participants are either matched based on all relevant characteristics explaining participation and outcomes, or;</p> <p>b) all relevant characteristics are accounted for.**</p> <p>c) and the data set used contains relevant variable that are measured in a relevant way (i.e. they were not collected for a different purpose initially and therefore are good proxy for some characteristics).</p> <p>**Accounting for and matching on all relevant characteristics is usually only feasible when the programme allocation rule is known and there are no errors of targeting. It is unlikely that studies not based on randomisation or regression discontinuity can score "YES" on this criterion. There are different ways in which</p>	<p>Score "Yes" if a) or b) and c) are satisfied</p> <p>Score "Probably yes" if a) or b) are addressed for but there is some doubt related to c), OR authors combined statistical matching and difference-in-difference to cope with unobservable differences, OR they only did statistical matching and there was clear rules for selection into the program (no self-selection).</p> <p>Score "Unclear" if · it is not clear whether all relevant characteristics (only relevant time varying characteristics in the case of panel data regressions) are controlled.</p> <p>Score "Probably no" if only a statistical matching was done and there was self-selection into the program.</p> <p>Score "No" if relevant characteristics are</p>

			<p>covariates can be taken into account. Differences across groups in observable characteristics can be taken into account as covariates in the framework of a regression analysis or can be assessed by testing equality of means between groups. Differences in unobservable characteristics can be taken into account through the use of instrumental variables (see also question 1.d) or proxy variables in the framework of a regression analysis, or using a fixed effects or difference-in-differences model if the only characteristics which are unobserved are time-invariant</p>	omitted from the analysis.
1: Selection bias - Justification	For identification based on an instrumental variable (IV estimation)	Open answer	<p>Score "Yes" if an appropriate instrumental variable is used which is exogenously generated: for example, due to a 'natural' experiment or random allocation.</p> <p>Score "Probably yes" if there is less evidence (no balance table showing differences between the intervention and comparison group).</p> <p>Score "Unclear" if the exogeneity of the instrument is unclear (both externally as well as why the variable should not enter by itself in the outcome equation).</p> <p>Score "Probably no" if there is evidence that enrolment in the program is correlated with a variable that might also have an effect on outcome and on the instrumental variable.</p>	

			Score "No" if it is clear that the instrument is not exogenous and affect the outcome through other channels than the program.	
2: Confounding - Assessment	2 - Group equivalence: was the method of analysis executed adequately to ensure comparability of groups throughout the study and prevent confounding?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear		
2: Confounding - Justification	For regression discontinuity design	Open answer	<p>a) The interval for selection of treatment and control group is reasonably small OR authors have weighted the matches on their distance to the cut-off point;</p> <p>b) and the mean of the covariates of the individuals immediately at both sides of the cut-off point (selected sample of participants and non-participants) are overall not statistically different based on t-test or ANOVA for equality of means;</p> <p>c) Significant differences in covariates of the individuals have been controlled in multivariate analysis; and for cluster-assignment, authors control for external cluster-level factors that might confound the impact of the programme.</p>	<p>Score "Yes, if criterion a), b), c) and d) are addressed.</p> <p>Score "Probably yes" if b) is not addressed but c) is addressed and differences in means are not large.</p> <p>Score "Unclear" if insufficient details are provided on controls; or if insufficient details are provided on cluster controls.</p> <p>Score "Probably no" if b) is not addressed (absence of a difference test or balance table) and there are doubt regarding the continuity on both sides of the cut-off point (a).</p> <p>Score "No" otherwise.</p>

2: Confounding - Justification	For non-randomised trials using difference-in-differences methods of analysis	Open answer	<p>a) The authors use a difference-in-differences (or fixed effects) multivariate estimation method;</p> <p>b) the authors control for a comprehensive set of individual time-varying characteristics, and for cluster-assignment, authors control for external cluster-level factors that might confound the impact of the programme**;</p> <p>c) and the attrition rate is sufficiently low and similar in treatment and control, or the study assesses that drop-outs are random draws from the sample (for example, by examining correlation with determinants of outcomes, in both treatment and comparison groups);</p> <p>**Knowing allocation rules for the programme – or even whether the non-participants were individuals that refused to participate in the programme, as opposed to individuals that were not given the opportunity to participate in the programme – can help in the assessment of whether the covariates accounted for in the regression capture all the relevant characteristics that explain differences between treatment and comparison</p>	<p>Score "Yes, if a, b, c, d (if relevant) are addressed and baseline imbalances between groups were relatively low OR the method was combined by a statistical matching.</p> <p>Score "Probably yes" if all possible variables are controlled for and the selection into the program was done according to clear rules, but baseline imbalances between groups were very large.</p> <p>Score "Unclear" if insufficient details are provided; or if insufficient details are provided on cluster controls.</p> <p>Score "Probably no" if some time-varying characteristics are not controlled for and the program was self-selected by the intervention groups.</p> <p>Score "No" if any of the criterion is not addressed.</p>
2: Confounding - Justification	<p>For statistical matching studies including propensity scores (PSM) and covariate matching**</p> <p>**Matching strategies are</p>	Open answer	<p>a) Matching is either on baseline characteristics or time-invariant characteristics which cannot be affected by participation in the programme; and the variables used to match are relevant (for example, demographic and socio-economic</p>	<p>Score "Yes, if a, b, c, and d (if relevant) are addressed.</p> <p>Score "Probably yes" if the selection into the program was done according to clear rules, which are used for the matching but</p>

	<p>sometimes complemented with difference-in-difference regression estimation methods. This combination approach is superior since it only uses in the estimation the common support region of the sample size, reducing the likelihood of existence of time-variant unobservable differences across groups affecting outcome of interest and removing biases arising from time-invariant unobservable characteristics.</p>		<p>factors) to explain both participation and the outcome (so that there can be no evident differences across groups in variables that might explain outcomes); and, for cluster-assignment, authors control for external cluster-level factors that might confound the impact of the programme b) in addition, for PSM Rosenbaum's test suggests the results are not sensitive to the existence of hidden bias; c) and, with the exception of Kernel matching, the means of the individual covariates are equated for treatment and comparison groups after matching; d) different matching methods including varying sample sizes yields the same results and authors take into account the use of control observations multiple times against the same treatment in their standard error calculation.</p>	<p>there are slight imbalances remaining after matching.</p> <p>Score "Unclear" if relevant variables are not included in the matching equation, or if matching is based on characteristics collected at end line; or if insufficient details are provided on cluster controls.</p> <p>Score "Probably no" if the program was self-selected by the intervention groups or participants OR if the selection into the program was done according to clear rules but there is no baseline data available to match the participants or groups on.</p> <p>Score "No" if matching was done based on variables that are likely to be affected by the program or any other scenario that affect a), b) c) or d).</p>
2: Confounding - Justification	<p>For regression-based studies using cross sectional data (excluding IV)</p>	Open answer	<p>a) The study controls for relevant confounders that may be correlated with both participation and explain outcomes (for example, demographic and socio-economic factors at individual and community level) using multivariate methods with appropriate proxies for unobservable covariates, and, for cluster-assignment, authors control particularly for external cluster-level factors that might confound the impact of the programme; b) and a Hausman test with an appropriate instrument suggests there is no evidence of endogeneity**;</p>	<p>Score "Yes, if a, b, c and d are addressed.</p> <p>Score "Probably yes" if all criterion are addressed but authors did not report the Hausman test (b).</p> <p>Score "Unclear" if relevant confounders are controlled but appropriate proxy variables or statistical tests are not reported; or if insufficient details are provided on cluster controls.</p> <p>Score "Probably no" if any of the criterion other than b) is not addressed.</p>

			<p>c) and none of the covariate controls can be affected by participation; d) and either, only those observations in the region of common support for participants and non-participants in terms of covariates are used, or the distributions of covariates are balanced for the entire sample population across groups;</p> <p>**The Hausman test explores endogeneity in the framework of regression by comparing whether the OLS and the IV approaches yield significantly different estimations. However, it plays a different role in the different methods of analysis. While in the OLS regression framework the Hausman test mainly explores endogeneity and therefore is related with the validity of the method, in IV approaches it explores whether the author has chosen the best available strategy for addressing causal attribution (since in the absence of endogeneity OLS yields more precise estimators) and therefore is more related with analysis reporting bias.</p>	<p>Score "No" if none of the criterion are addressed.</p>
--	--	--	--	---

2: Confounding - Justification	For identification based on an instrumental variable (IV estimation)	Open answer	<p>a) The instrumenting equation is significant at the level of $F \geq 10$ (or if an F test is not reported, the authors report and assess whether the R-squared (goodness of fit) of the participation equation is sufficient for appropriate identification);</p> <p>b) the identifying instruments are individually significant ($p \leq 0.01$); for Heckman models, the identifiers are reported and significant ($p \leq 0.05$);</p> <p>c) where at least two instruments are used, the authors report on an over-identifying test ($p \leq 0.05$ is required to reject the null hypothesis); and none of the covariate controls can be affected by participation and the study convincingly assesses qualitatively why the instrument only affects the outcome via participation. If the instrument is the random assignment of the treatment, the reviewer should also assess the quality and success of the randomisation procedure in part a).</p> <p>d) and, for cluster-assignment, authors particularly control for external cluster-level factors that might confound the impact of the programme (for example, weather, infrastructure, community fixed effects, and so forth) through multivariate analysis.</p>	<p>Score "Yes, if a, b, c, d (if relevant) are addressed.</p> <p>Score "Probably yes" if one of the test required for criterion a) or b) is not reported but the other is, and the rest of the criterion are addressed and the instrument is convincing.</p> <p>Score "UNCLEAR" if relevant confounders are controlled for but appropriate statistical tests are not reported; or if insufficient details are provided on cluster controls</p> <p>Score "Probably no" if exogeneity of the instrument is not convincing and appropriate tests are not reported.</p> <p>Score "No" otherwise if any of the tests required for criterion a), b) or c) are reported and not satisfied.</p>
-----------------------------------	--	-------------	--	---

3: Performance bias - Assessment	3 - Performance bias: was the process of being observed free from motivation bias?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	<p>a) For data collected in the context of a particular intervention trial (randomised or non-randomised assignment), the authors state explicitly that the process of monitoring the intervention and outcome measurement is blinded, or argue convincingly why it is not likely that being monitored could affect the performance of participants in treatment and comparison groups in different ways (such as resulting in Hawthorne or John Henry effects).</p> <p>b) The study is based on data collected in the context of a survey, and not associated with a particular intervention trial, or data are collected from administrative records or in the context of a retrospective (ex post) evaluation.</p>	<p>Score "Yes" if either criterion a) or b) are satisfied;</p> <p>Score "Probably yes" if the study is based on survey data collected during a trial and there is no obvious issue with the monitoring processes but authors do not mention potential risks.</p> <p>Score "Unclear" if it is not clear whether the authors use an appropriate method to prevent Hawthorne and John Henry Effects (e.g. blinding of outcomes and, or enumerators, other methods to ensure consistent monitoring across groups). Hawthorne effects may result where participants know that they are being observed and John Henry Effects may result from participant knowledge of being compared.</p> <p>Score "Probably no" if there was imbalance in the frequency of monitoring in intervention groups, which might have influenced participants' behaviours.</p> <p>Score "No"</p>
3: Performance bias - Justification	Performance bias - Justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	

<p>4: Spill-overs, cross-overs and contamination - Assessment</p>	<p>4 - Spill-overs, cross-overs and contamination: was the study adequately protected against spill-overs, cross-overs and contamination?</p>	<p>1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear</p>	<p>a) There was no implementation issues that might have led the control participants to receive the treatment (implementer's mistake). b) The intervention is unlikely to spill-over to comparisons (e.g. participants and non-participants are geographically and/or socially separated from one another and general equilibrium effects are not likely) or the potential effects of spill overs were measured (e.g. variation in the % of unit within a cluster receiving the treatment). c) There is no risk of contamination by external programs: the treatment and comparisons are isolated from other interventions which might explain changes in outcomes. d) There is nothing in the surveys that might have given the control participants an idea of what the other group might receive OR they did but there is no risk that this has changed their behaviours; AND the survey process did not reveal information to the control group that they did not have before (e.g. the study aims to measure increase in take up of a service or product that participants might not know about) Authors might put something in place in the design of the study that allows to control for that survey effect (e.g. a pure control with no monitoring except baseline end line)</p>	<p>Score "Yes" if criterion a), b), c) and d) are satisfied; Score "Probably yes" if there is no obvious problem but there is no information reported on potential risks related to spill overs, contamination, or survey effects in the control group OR if there were issues with spill-overs but they were controlled for or measured. Score "Unclear" if spill-overs, cross-overs, survey effects and/or contamination are not addressed clearly. Score "Probably no" if any of the criterion a), b), c) or d) are not satisfied but the scale of the issue is not clear. Score "No" if any of the criterion a), b), c) or d) are not satisfied and happened at a large scale in the study.</p>
---	---	---	--	---

4: Spill-overs, cross-overs and contamination - Justification	Spill-overs, cross-overs and contamination - Justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	
5: Outcome measurement bias - Assessment	5 - Outcome measurement bias	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	<p>a) Outcome assessors are blinded or the outcome measures are not likely to be biased by their judgement.</p> <p>b) For self-reported outcomes: respondents in the intervention group are not more likely to have accurate answers due to recall bias;</p> <p>c) For self-reported outcomes: respondents do not have incentives to over/under report something related to their performance or actions, OR researchers put in place mechanisms to reduce the risk of reporting bias (researchers not strongly involved in the implementation of the program and it is clear that their answers to the survey will not affect what they receive in the future) OR authors have measured the risks of bias through falsification tests or measuring the effect on placebo outcomes in cases where there was a risk of reporting bias.</p> <p>d) Timing issue: the data collection period did not differ between intervention and comparison group, the baseline data is not likely to be affected by the beginning of the intervention or affects a small percentage of the study participants.</p>	<p>Score "Yes" if criterion a), b), c) and d) are satisfied:</p> <p>Score "Probably yes" if there is a small risk related to any of a), b), c) or d) and there is no more information provided to justify the absence of bias OR if there was a high risk of bias but authors have either controlled it in their design or measured it with a placebo outcomes.</p> <p>Score "Unclear" if it there is a high risk related to any of a), b), c) or d) and there is no more information provided to justify the absence of bias.</p> <p>Score "Probably no" if there are high risk related to a), b), c) or d) and it is clear that authors were not able to control for this bias.</p> <p>Score "No" if there is evidence of bias.</p>

5: Outcome measurement bias - Justification	Outcome measurement bias - Justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	
6: Reporting bias - Assessment	6 - Selective analysis reporting: was the study free from selective analysis reporting?	1= Yes, 2 = Probably Yes, 3 = Probably No, 4 = No, 8 = Unclear	<p>a) a pre-analysis plan is published, especially for prospective NRS but it should also be for retrospective studies</p> <p>b) authors use 'common' methods of estimation (i.e. credible analysis method to deal with attribution given the data available) ;</p> <p>c) There is no evidence that outcomes were selectively reported (e.g. results for all relevant outcomes in the methods section are reported in the results section) ;</p> <p>d) Requirements for specific methods of analysis:</p> <ul style="list-style-type: none"> - For PSM and covariate matching: (a) Where over 10% of participants fail to be matched, sensitivity analysis is used to re-estimate results using different matching methods (Kernel Matching techniques); (b) For matching with replacement, no single observation in the control group is matched with a large number of observations in the treatment group. - For IV (including Heckman) models, (a) The authors test and report the results of a Hausman test for exogeneity ($p \leq 0.05$ is required to reject the null hypothesis of exogeneity); (b) the coefficient of the selectivity correction term (ρ) is significantly different from zero ($P < 0.05$) (Heckman approach). 	<p>Score "Yes" if a), b), c) and d) are satisfied OR if a) is not met and it is a retrospective NRS.</p> <p>Score "Probably Yes" if authors combined methods and reported relevant tests (d) only for one method OR if all the criteria are met except for a) and it is a prospective NRS</p> <p>Score "Unclear" if intended outcomes not specified in the paper OR if any of the requirements for d) are not reported.</p> <p>Score "Probably No" if b) is addressed, but authors did not present results for all outcomes announced in the method section OR did not meet requirement d) although reported.</p> <p>Score "No" if authors use uncommon or less rigorous estimation methods such as failure to conduct multivariate analysis for outcomes equations OR if some important outcomes are subsequently omitted from the results or the significance and magnitude of important outcomes was not assessed.</p>

			- For studies using multivariate regression analysis, authors conduct appropriate specification tests (e.g. testing robustness of results to the inclusion of additional variables, or (very rare) reporting results of multicollinearity test etc).	
6: Reporting bias - Justification	Analysis reporting bias - Justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	
7: Other bias - Assessment	7 - Other risks of bias: Is the study free from other sources of bias?	1= Yes, 4 = No	Score "Yes" if the reported results do not suggest any other sources of bias. Score "No" if other potential threats to validity are present, and note these here (e.g. coherence of results, survey instruments used are not reported)	
7: Other bias - Justification	Other risks of bias - Justification	Open answer	Justification for coding decision (Include a brief summary of justification for rating, mentioning your response to all sub questions, cite relevant pages).	
8: External validity	8 - External validity	Open answer	Open answer- what do authors say about external validity, if anything?	