



UNIVERSITÀ DEGLI STUDI DI PADOVA

DIPARTIMENTO DI PSICOLOGIA DELLO SVILUPPO E DELLA SOCIALIZZAZIONE

CORSO DI LAUREA MAGISTRALE IN PSICOLOGIA CLINICO-DINAMICA

TESI DI LAUREA MAGISTRALE

**Sample size planning informed by experts and
literature: toward clinically meaningful effect sizes.**

Pianificazione della numerosità campionaria basata sul giudizio degli esperti e sintesi della letteratura. Una proposta per stimare dimensioni dell'effetto clinicamente rilevanti.

Relatore:

PROF. GIANMARCO ALTOÈ

Laureando:

EMANUELE BOLLINI

MATRICOLA: 2119954

Correlatori:

DOTT.SSA AMBRA PERUGINI

Anno Accademico 2024/2025

Acknowledgements

I began this work believing the acknowledgements would be a brief formality, but now I find myself in the uncomfortable position of feeling unable to convey my deep gratitude with a couple of paragraphs.

I am sincerely grateful to Professor Cristea for her central role in establishing the foundational conditions for this research.

Also, my most sincere thanks go to those who directly supervised this work.

To Professor Altoè, for demanding such rigour with such kindness. I have deeply enjoyed learning from your expertise.

To Ambra for your reassuring and insightful contributions that made all the difference, from brainstorming phases to refining the details.

To Laura, for your meticulous commitment that went far beyond any formal obligation.

To Matteo, For your unexpected availability, which was a valuable help for the final stages.

Finally, my gratitude also extends to all the other, indirect supporters of this work, whom I will thank elsewhere.

Table of contents

1	Replication crisis in psychology	5
1.1	Introducing the replication crisis	5
1.2	What is a replication	6
1.3	Replication crisis	7
1.4	A broad conceptual crisis	8
1.4.1	Validity Crisis	8
1.4.2	Theory crisis	9
1.5	Questionable Research Practices	10
1.6	Inferential framework and its misunderstanding	11
1.6.1	Frequentist paradigm	11
1.6.2	Neyman and Pearson approach	12
1.6.3	<i>NHST</i> approach	13
1.7	Aim of this thesis	15
2	Enhance effect size interpretation	17
2.1	Why we need effect sizes	18

2.2	Effect size	18
2.2.1	Illustrating effect size	19
2.3	Interpretation of effect sizes	21
2.3.1	Statistical significance and practical significance	21
2.3.2	The misuse of Cohen’s classification benchmarks	22
2.4	Alternative approaches to interpret effect sizes	24
2.4.1	The plausible effect size	25
2.4.2	The smallest effect size of interest	26
2.5	A methodology for building a plausible <i>SESOI</i>	27
2.6	Conclusions	31
3	Expert Elicitation	33
3.1	Why we need expert elicitation	34
3.2	What is expert elicitation	35
3.2.1	The matter of expertise	36
3.2.2	The matter of estimates	36
3.2.3	The aggregation problem	37
3.2.4	The matter of biases	38
3.3	The leading protocols	39
3.4	Why eliciting a <i>SESOI</i>	42
3.5	A tailored procedure for eliciting <i>SESOI</i>	42
4	Methodology	44
4.1	Introduction	44

4.2	Methods and materials	46
4.2.1	Meta-analysis	47
4.2.2	Elicitation procedure	47
4.2.3	Comparing <i>SESOI</i> and Plausible effect	52
4.3	Analysis	57
4.4	Results	57
4.5	Discussion	58
4.5.1	Conclusion	60
5	Conclusions	61
5.1	General results	62
5.2	Limitations	63
5.2.1	Elicitation related	64
5.2.2	Research context related	65
5.2.3	Statistical limitations	66
5.3	Future developments	66
	References	68
	Appendices	78
A		78
B		81

Chapter 1

Replication crisis in psychology

In this chapter we outline key methodological and conceptual issues underlying the replication crisis, introducing the concept of replication and its role in science, then addressing major problems such as low replication rates, weak theories, poor measurement, and misleading statistical practices. After delimiting the inferential paradigm, we will then focus on the misuse of the inferential paradigm, especially the overreliance on Null Hypothesis Significance Testing (*NHST*).

1.1 Introducing the replication crisis

In recent years, psychology has faced a credibility crisis, raising doubts about the credibility of many published findings (Ioannidis, 2005). Since the beginning of such crisis, significant efforts have been devoted to improving data analysis, theory testing, and Open Practices (Lakens, 2019; Nosek et al., 2015; Nosek & Lakens, 2014; Oberauer & Lewandowsky, 2019).

We align with the interventionist perspective that identifies superficial methodological choices (specifically, the neglect of rigorous sample size planning and clinically meaningful effect size inter-

pretation) as one of the main causes of the crisis. In this thesis, we will develop and substantiate this explanation, arguing that a principled application of the Neyman-Pearson framework provides the necessary corrective to these methodological shortcomings.

1.2 What is a replication

“Scientific progress is a cumulative process of uncertainty reduction that can only succeed if science itself remains the greatest skeptic of its explanatory claims” (Abraham et al., 2015, p. 7).

Arguably, the very foundation of science depends on the replication of studies (Murphy et al., 2025). The greater our ability to replicate findings over time, the more confidence we gain in our beliefs (Collaboration, 2015; Schmidt, 2009). Therefore, replication allows us to be more skeptical about certain theories (Eronen & Bringmann, 2021), while also more confident in confirming and generalizing others (Cohen, 1994; Schmidt, 2009).

It is important to define terms that are often used interchangeably in the literature, despite referring to distinct concepts. A *reproduction* study consists in observing the same results, performing the same statistical analysis with the same data (indeed, it’s also often referred to as “computational reproducibility”). In contrast, *replication* study consists in conducting the same study with different data and obtaining similar results (Nosek et al., 2022). The main challenge in replication studies is to determine whether the results we obtain are sufficiently close to the original ones, to consider the replication successful.

Replication is often divided in two categories. Replication that aims to closely match the original study is called *direct replication*. Direct replication are very useful to confidently identify false positive (Nosek & Lakens, 2016). On the contrary, when variations are introduced in the study (e.g., different population), this is referred to as *conceptual replication* (Derksen & Morawski,

2022). Conceptual replication may be useful to explore if the phenomenon can occur in different scenarios, thereby enhancing our comprehension of the effect, and offering more support to the theory (Nosek & Lakens, 2016). However, some argue that conceptual replications often attribute the effect to changes in study design rather than challenging the phenomenon itself (Nosek & Lakens, 2016). Offering a different perspective on the matter, Machery (2020) proposes abandoning the direct/conceptual distinction, and proposes instead to classify replications on how many components are resampled (e.g., units, treatments, measurements, settings): each combination of resampled components examines a different aspect of the original experiment’s reliability.

1.3 Replication crisis

In recent years, many studies have shown low replication rates in psychological research. This has led to the conclusion that many psychological findings may actually be false (Collaboration, 2015; Ioannidis, 2005; Machery, 2020; Malich & Munafò, 2022; Oberauer & Lewandowsky, 2019; Vazire et al., 2022).

The main reason why positive results can be easily produced, and therefore published, is due to some inappropriate methodological practices: the most relevant among these are the misuse of an inferential procedure known as the Null Hypothesis Significance Test (*NHST*), and a range of questionable research practices (QRPs). A larger overview regarding of these two phenomena will be provided in the dedicated section (Head et al., 2015; Ioannidis, 2005; Scheel, 2022).

The spread of false positives is worsened by a publication bias. Publications offering significant results tend to be published more frequently (Scheel, Schijen, et al., 2021). There’s a substantial pressure to publish, and studies presenting significant effects are often published at higher rate (Banks et al., 2016; Primbs et al., 2023). At the same time, likely due to the inherent difficulty

in replicating studies in psychology (Hall III, 2023), scientific publishing world prioritizes novelty over replication (Collaboration, 2015). Which contributes to the dissemination of low-quality or non-replicable findings.

This controversial finding undermines the very foundation of science, as replication is typically the primary mean by which we identify theories that are likely true. However, since replication is rarely done in psychology, overlapping theories that have not been decisively falsified yet remain common in the field (Eronen & Bringmann, 2021).

1.4 A broad conceptual crisis

The replication crisis has drawn attention to additional issues that characterize the field of psychology, many of which originate from the discipline’s conceptual foundations and are worth noting: the validity crisis and the theory crisis. These aspects are broad and have been extensively discussed elsewhere. Here, a brief explanation will be provided, for the sake of completeness, before proceeding with the discussion of the replication crisis.

1.4.1 Validity Crisis

Validity crisis concerns the precise definition of the object of study within a given field (Eronen & Bringmann, 2021). If researchers are unable to define what they aim to measure, it will be unclear whether they will have been measuring the intended construct in the first place (Flake & Fried, 2020). In recent years, in response to the replication crisis, increasing attention has been given to the foundations of psychological measurement, as many studies fail to provide clear definitions of the constructs of interest (Vazire et al., 2022). Earlier, Flake and Fried (Flake & Fried, 2020) introduced the concept of *questionable measurement practices* (QMP), addressing all those practices

that cast doubt on the validity of measurement, such as lack of transparency in how measures were used and superficial or inappropriate application of measurement tools.

Following the recognition of this crisis, several recommendations have been proposed to improve measurement quality. For instance, Flake (Flake & Fried, 2020) emphasizes the need for greater transparency, and provides a structured series of questions to help researchers define construct and avoid QMP.

The crisis of measurement validity is also reflected in the weakness of psychological theories, which are often based on poorly defined constructs.

1.4.2 Theory crisis

Theory crisis refers to “a weak logical link between theories and their empirical test” (Oberauer & Lewandowsky, 2019, p. 2). This is arguably a more fundamental problem: theories are not strong enough to derive strong hypotheses or to predict the dimension of any effect (Eronen & Bringmann, 2021; Fried, 2020). Without strong hypotheses, theories cannot be falsified. Moreover weak theories can easily be defended after the results are known by adding auxiliary hypotheses. A common line of defense is the appeal to “hidden moderators”, which are often invoked to explain failed replication (Fried, 2020). If that is the case though, we are not able to tell whether data corroborates our theory or not.

As Eronen and Bringmann highlighted we should invest more energy into building stronger and more precise theories (Eronen & Bringmann, 2021). However psychology is a challenging field, and to build strong theories we first need to identify clear phenomena. Therefore, some argue they should focus more on exploratory research (Fife & Rodgers, 2022; Oberauer & Lewandowsky, 2019; Scheel, Tiokhin, et al., 2021), and in phenomenon-driven research, as it helps constrain the space

of plausible theories (Eronen & Bringmann, 2021).

While the validity crisis and the theory crisis certainly warrant further investigation to strengthen psychology’s conceptual foundations, the replication crisis is more commonly attributed to other factors. The two main areas of concern are questionable research practices and misunderstandings related to inferential methodology.

1.5 Questionable Research Practices

Questionable Research Practices (QRPs). QRPs are all the methods implemented in the field that aim to obtain a significant value, such as changing the hypothesis after the results are known (often referred to as “HARKing”) (Kerr, 1998), stopping data collection after achieving the desired result (often referred to as “optional stopping”), collecting more data after seeing whether results were significant, manipulating statistical analysis in order to obtain significance (often referred to as “ p -hacking”), presenting the results of a study that best support the hypothesis instead of reporting all the findings (oftentimes also called “cherry picking”) and many others (John et al., 2012; Scheel, Schijen, et al., 2021). Therefore, QRPs include unethical and ambiguous practices, and have been shown to increase the likelihood of obtaining significant results (John et al., 2012), and increase the spread of false positive (Flake & Fried, 2020). Although often not driven by deliberate intent, QRPs are typically fueled by self-serving biases and the strong pressure to produce significant findings (John et al., 2012; Simmons et al., 2011).

To improve the situation, several solutions have been recently proposed, including strengthening statistical standards, conducting high-powered replications, providing open data, materials and algorithms to enhance reproducibility, and clearly distinguishing a priori hypothesis through preregistration (a practice aimed at preventing HARKing and clarifying the distinction between

confirmatory and exploratory research) (Nosek et al., 2015; Oberauer & Lewandowsky, 2019)

Significant results, however, refer primarily to outcomes derived from the most widely used statistical framework in psychology: null hypothesis significance testing (*NHST*) (Head et al., 2015). A brief clarification is useful here to understand its role in the replication crisis.

1.6 Inferential framework and its misunderstanding

To address the methodological issues at the core of the replication crisis, we will first specify the inferential framework we are adopting. Then, we will outline the Neyman-Pearson approach, as it is useful for understanding the weaknesses of the aforementioned *NHST*.

1.6.1 Frequentist paradigm

As previously noted, we must first clarify the logic used to interpret empirical evidence. Among all the approaches used in inferential sciences, two in particular are mostly used, namely frequentist approach and bayesian approach (Dienes, 2008; Van de Schoot et al., 2014). Although this thesis works within a frequentist framework, both approaches are sound; see Van de Schoot et al. (2014) for an overview of the Bayesian perspective.

In a frequentist perspective, probability is defined as the frequency (i.e., number of occurrences) of an event in a set period of time (VandenBos, 2007). Having set the inferential paradigm of reference, we will now outline the problems that have emerged in the context of hypothesis testing, and which solutions have been proposed.

1.6.2 Neyman and Pearson approach

Within this inferential framework, several methodologies exist for hypothesis testing. One of the most studied is the Neyman-Pearson approach.

In the Neyman-Pearson approach, before collecting data, a predetermined level of inferential risk is established in order to make a decision between different hypotheses defined a priori. The test result will allow one to act as if one of the two hypotheses were true and the other false, given a certain level of risk. To do this, one must decide in advance on the following elements (Gigerenzer, 2004):

- 1) Defining two opposing hypothesis. The null hypothesis (H_0), which usually assumes no effect or difference, and the alternative hypothesis (H_1), which posits the presence of an effect and is the focus of testing. The magnitude of the effect has to be decided basing on theoretical or practical criteria. By defining H_0 and H_1 , the associated sampling distributions are also specified. This allows for the division of the sample space into acceptance and rejection regions for each hypothesis.
- 2) To define the risks associated with the test. To do so, we specify in advance values for α (the probability of incorrectly rejecting the null hypothesis when it is true, namely Type I error) and β (the probability of failing to reject the null hypothesis when it's false). Power can be then calculated as the probability that the test has to reject the null hypothesis (H_0) when the alternative hypothesis (H_1) is true (Altoè et al., 2020). The level of risk one intends to accept depends on the context and the phenomenon under investigation (Maier & Lakens, 2022).
- 3) Define the sample size. Once we know the a priori effect size, and we determine the values for α and β , the sample size is consequently determined in a design analysis.

The test then yields a p -value, which is the probability of observing the same or higher data than observed ones, if H_0 is true. The p -value is then compared against the pre-specified α (therefore also called “significance level”), and we either accept H_1 and reject H_0 , or reject H_1 and accept H_0 (Neyman, 1957).

The Neyman and Pearson actually originates from an other approach, the Fisher approach (Fisher, 1955). This approach only yields the p -value under H_0 , so the smaller the p -value, the more the data are unlikely under the Null Hypothesis. This approach though is intended for very preliminary analyses, and results must be interpreted basing on the context at hand (Gigerenzer, 2004). Specifying the usage and limitation of these procedures should make it easy to clarify the issues related to the previously noted *NHST*, which we will outline in the following section.

1.6.3 *NHST* approach

During its history, the development of psychology as a field has led to the widespread use of a simplified and less rigorous variation of this method, the aforementioned Null Hypothesis Significance Testing (*NHST*) (Gigerenzer, 2004). This test represents a hybrid form of the Neyman-Pearson approach and the earlier approach developed by Fisher (Fisher, 1955). The *NHST* paradigm works as follows (Gigerenzer, 2004):

- H_0 is defined as a null hypothesis (e.g., we test against the complete absence of an effect).
- No alternative hypothesis is specified.
- α is conventionally and uncritically set at 0.05.
- Finally, if $p\text{-value} < 0.05$ (statistically significant), the null hypothesis is rejected.

As previously noted, this methodology constitutes a simplified and weaker version of the Neyman-Pearson framework, as it omits two of its essential components. First, the significance level (α) is

commonly fixed at 5% without theoretical justification, making it an uncritical convention rather than a reasoned decision. In fact, it may be more important to identify false negatives more precisely in different context, e.g. a false positive that leads to a useless surgery may arguably be worse than uselessly having some therapy sessions, where we may accept higher risk for false positive. Second, since no alternative hypothesis is specified, no hypothesis can ultimately be accepted. This absence also prevents the definition of an ideal sample size in advance, often leading to under-powered studies (Gigerenzer, 2004). Moreover, without an alternative hypothesis, the a priori dimension of the effect cannot be calculated, making the interpretation more misleading.

To these characteristics, one must add a series of further elements of concern, as the *NHST* approach is characterized by a widespread misconception about logic of hypothesis testing and interpretation of p -values (Fife & Rodgers, 2022; Lakens, 2021). See Lakens (2016) for a more detailed discussion.

A vast majority of misconceptions consists in attributing to p -values informative strength over the theory to be tested (Gigerenzer, 2004). In fact, once the test yields a significant result, the research hypothesis, which had not been previously specified, is usually considered confirmed. This is incorrect, as no alternative theory had been previously specified. On the other hand, if a non-significant result is obtained, most psychologists often conclude that there is no effect, which is equally false (Lakens, 2017). Truly, a significant result does not indicate that a theory is likely true; it only tells us how likely the observed data are, assuming the null hypothesis is true (Cohen, 1994), since, as Gigerenzer states, “The probability $P(D|H_0)$ is not the same as $P(H_0|D)$, and more generally, a significance test does not provide a probability for a hypothesis” (2004, p.95). Ultimately, p -value is often erroneously considered as the strength of an effect or relationship (Head et al., 2015).

In synthesis, even when properly understood, *NHST* presents some unavoidable weaknesses. It only allows for the rejection of the null hypothesis, offering little to no meaningful information about the interpretation of the effect (Cohen, 1994). Moreover, since population parameters in the real world are virtually never exactly zero, even negligible differences can become statistically significant with sufficiently large sample sizes (Lin et al., 2013). This can lead to a misleading sense of confirmation for effects that are practically irrelevant (Lin et al., 2013; Malich & Munafò, 2022). Additionally, by omitting the alternative hypothesis and the corresponding effect size, *NHST* prevents the calculation of an appropriate sample size to achieve a desired level of statistical power (Gigerenzer, 2004).

1.7 Aim of this thesis

As shown in this chapter, improving inferential methodology and addressing the replication crisis requires a broad set of interventions. One intervention is strengthening inference by applying the Neyman-Pearson framework as originally intended. Doing so, this thesis directly addresses two critical aspects of robust inference, often neglected in recent times: rigorous sample size planning and meaningful effect sizes interpretation.

To do this effectively, it is essential to pre-specify and formalize both the null and alternative hypotheses when designing a statistical test (as we outlined, when one is not able to formalize both hypothesis, more exploratory research should be done.) This approach avoids some of the major limitations of Null Hypothesis Significance Testing (*NHST*). As discussed, *NHST* does not support prospective power analysis, which is necessary to determine an adequate sample size. As a result, many studies are underpowered (Cohen, 1962). Moreover, statistical significance is sometimes achieved for trivially small effects, leading to results that are only weakly informative. Finally,

researchers frequently misinterpret p-values as indicators of the strength of an effect.

As we have shown, many of these issues can be addressed by explicitly considering the effect size under the alternative hypothesis. This allows for proper planning of adequately powered studies and requires a theoretical estimation of the expected effect, thereby emphasizing the importance of effect sizes.

This thesis addresses the central challenge of specifying a meaningful alternative hypothesis by proposing a possible solution to facilitate its pre-specification.

In chapter two, we will demonstrate that interpreting effect sizes is crucial for obtaining meaningful results, although they are often assessed using uninformative benchmarks. We will introduce two approaches to ensure more meaningful consideration of effect sizes. The first is the “plausible effect size” (Gelman & Carlin, 2014), which helps define an effect size based on what is realistically expected in the context of the specific test. The second is the “Smallest Effect Size of Interest” (Anvari & Lakens, 2021), also known as *SESOI*, a concept recently applied in psychology that can help identify practically relevant effects. Finally, we will propose a method to compare these two approaches, assessing whether the effect sizes we consider meaningful are likely to be detected in practice.

Having established the foundation of this thesis, in chapter three we will explore expert elicitation procedures, methodologies for deriving effect sizes based on expert judgment, which we will use to determine the *SESOI* for the present research.

Finally, in chapter four, we bridge theory and practice by implementing the proposed framework. We will ultimately demonstrate its practical utility by applying the integrated results to inform a design analysis.

Chapter 2

Enhance effect size interpretation

"Thinking hard about effect sizes is important for any school of statistical inferences [i.e., Frequentist or Bayesian], but sadly a process often neglected."

— Dienes, 2008, p.92

In this chapter we will examine the effect size, providing a definition and a brief explanatory example. We will then outline its interpretation challenges, addressing the limits of statistical significance compared to the concept of practical significance. The distinction will allow us to introduce methods for formulating meaningful hypotheses, therefore explaining the plausible effect size and the smallest effect size of interest (*SESOI*) to enhance inferential robustness. Then we will propose a practical methodology as a combination of the two above.

2.1 Why we need effect sizes

In the previous chapter, we highlighted the importance of formalizing an alternative hypothesis while doing hypothesis testing within the Neyman-Pearson framework, or, in other words, the importance of specifying an effect size for the study. As outlined, this would improve research in the social sciences in several ways. Firstly, by allowing us to conduct a proper design analysis, it would enable transparent preplanning of an appropriate sample size, thus reducing the risk of false positives (Altoè et al., 2020). Secondly, the ability to predict effects and then observe whether they are confirmed or disconfirmed would lead to more falsifiable studies (Anvari & Lakens, 2021). Furthermore, since p -values alone can lead to the acceptance of results with low practical significance, incorporating effect sizes reduces the risk of attributing importance to potentially meaningless findings. These advantages highlight why effect sizes are essential complements to significance testing in study design and in result interpretation.

2.2 Effect size

As Pek and Flora report (2018, p.209), “the term *effect size* literally translates to some magnitude (or size) of the impact (or the effect) of a predictor or an outcome variable”. Briefly said, it quantifies a given effect. As the definition depends on the phenomenon of interest, there is a very wide amount of possible definition to determine an effect size (Kelley & Preacher, 2012). In psychology research, though, effect sizes are most often used to indicate a relationship or a difference between two variables (Borenstein et al., 2021).

There are two main ways to report effect sizes in the literature. Firstly, the raw mean difference, where the effect size is expressed in the same units as the original measurement scale

(Borenstein et al., 2021). For example, the raw mean difference could quantify the mean difference between two groups on blood pressure measurement. Raw effect sizes offer a direct and intuitive idea of the meaning of the effect, especially given a widespread unit scale (as mmHg for the given example)(Borenstein et al., 2021). Nevertheless, many researches in social sciences use a variety of different measurements, adopting different scales, making it challenging to have a shared understanding of raw effect sizes (Altoè et al., 2020).

Given this limitation, researchers often prefer a second approach, i.e. standardized effect sizes, which facilitate comparison across studies and contexts. For this purpose, many different standardized effect sizes exist. In addressing this matter in the present thesis, however, we will focus on the most widely used and well-known measure (Altoè et al., 2020), namely Cohen’s d (μ_d) (Cohen, 1988). In a concise formulation, Cohen’s d is the raw difference between two population means ($\mu_A - \mu_B$) divided by the common standard deviation (σ):

$$\delta = \frac{\mu_A - \mu_B}{\sigma}$$

This formulation allows for standard interpretation across different measurement scales, which enhances comparability in psychological research.

To illustrate how a standardized effect size such as Cohen’s d operates in practice, we provide a practical illustration in the following section.

2.2.1 Illustrating effect size

To illustrate how a standardized effect size such as Cohen’s d operates in practice, we now present a simplified example based on a pre-post intervention study in a clinical context. In such studies, psychological distress is typically measured before and after an intervention using standardized in-

struments, for example, the Kessler Psychological Distress Scale (K-10). Participants are randomly assigned to a treatment group or a control group, and the effect of the intervention is assessed by comparing changes in K-10 scores. In the following example, intended solely for didactic purposes, we provide a visual and conceptual representation of what it means to specify an effect size under the null and alternative hypotheses. Specifically, our alternative hypothesis specifies an expected standardized effect of $d = 0.41$.

The figure 2.1 shows two theoretical distributions of the test statistic under the null hypotheses (H_0) and alternative hypotheses (H_1), within a Neyman-Pearson framework. These distributions derive from a comparison between a treatment group and a control group in a pre-post design. The black curve represents the distribution under the null hypothesis (H_0), which assumes that the intervention has no effect (i.e., the mean change in the treatment group is equal to that in the control group) corresponding to an effect size of $d = 0$. The blue curve represents the alternative hypothesis (H_1), assuming a true standardized effect size of $d = 0.41$, meaning the treatment group improves, on average, by 0.41 standard deviations more than the control group. The distance between the two distributions reflects the assumed effect magnitude.

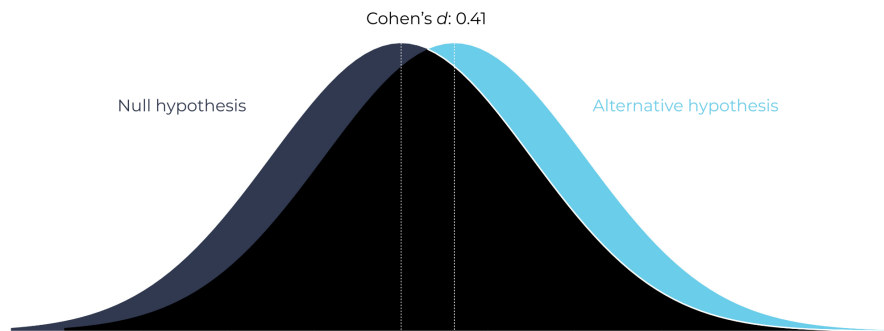


Figure 2.1: A graphical representation of a Cohen's d , in a RCT group (Magnusson, 2023).

We have provided a graphic illustration of an effect size, which represents values distribution when the treatment has no effect (H_0) and when the effect is 0.41 (H_1). It is quite intuitive though that neither a graphical representation nor a numerical value clearly informs us of *how much* of improvement patient have experienced under the alternative hypothesis (Funder & Ozer, 2020). This example introduces the main matter discussed in this thesis that we will detail in the following section, that is interpretation of effect size.

2.3 Interpretation of effect sizes

When addressing effect sizes, it is important to distinguish between informative and uninformative ones. While this may not be necessary in all sciences, it goes without saying that in psychology we must be able to differentiate the magnitude of effects. Taking into account the clinical context, one in which this distinction is especially critical, it is clear that we need to discriminate between mild and highly effective interventions, just as we distinguish between severe and modest clinical cases.

However, in psychology, since there is no guidance for interpreting effect sizes (Ferguson, 2016), there is a tendency to treat any effect size as a success (Anvari et al., 2023; Hilgard, 2021), which hinders a clear distinction between these two categories. To counter this issue, it is useful in this regard introducing the distinction between practical and statistical significance, which are not alternatives, but complementary dimensions of result interpretation.

2.3.1 Statistical significance and practical significance

Statistical significance is met when a certain p -value indicates that the result meets the level of evidence required by the researcher's chosen threshold (α) to reject the null hypothesis (H_0), but p -value alone provides no information about the magnitude or practical relevance of the observed

effect (Boring, 1919; Kelley & Preacher, 2012; Pek & Flora, 2018). A statistically significant result may correspond to a trivially small effect, especially in studies with large sample sizes, while relevant effects may go undetected in underpowered studies. Practical significance, instead, refers to the distinction between interesting and uninteresting effects (Anvari & Lakens, 2021). In fact, as already shown, effect sizes themselves are not automatically informative of anything, and should be interpreted in light of substantive criteria (Grissom & Kim, 2012; Primbs et al., 2023).

2.3.2 The misuse of Cohen’s classification benchmarks

These points illustrate the necessity of evaluating effect sizes within the specific context and goals of research. Despite this, the interpretation of effect sizes is usually related to a classification proposed by Cohen (Cohen, 1988), who suggested values of $d = 0.2$, $d = 0.5$, $d = 0.8$ as indicative of “small”, “medium” and “large” effect sizes respectively.

A graphical representation of the three effect sizes is reported at figure 2.2, 2.3 and 2.4 respectively. These sets of distributions derive from a hypothetical comparison between a treatment group and a control group in a pre-post design. The black curves always represent the distribution under the null hypothesis (H_0), which assumes that the intervention has no effect (i.e., the mean change in the treatment group is equal to that in the control group) corresponding to an effect size of $d = 0$. The blue curves always represent the alternative hypothesis (H_1), assuming a true standardized effect size of $d = 0.2$, $d = 0.5$, $d = 0.8$, meaning the treatment group improves, on average, by 0.2, 0.5 and 0.8 standard deviations more than the control group, respectively. The distance between each pair of distributions reflects the assumed effect magnitudes.

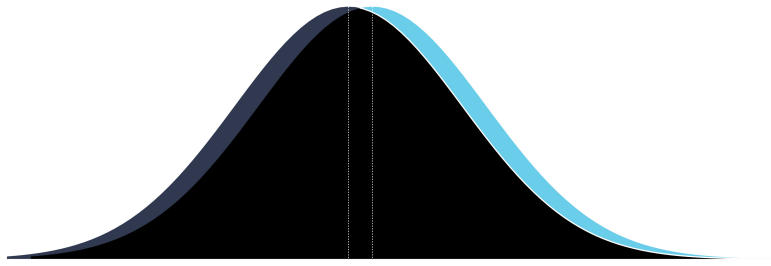


Figure 2.2: A graphical representation of 0.2 Cohen's d (Magnusson, 2023)

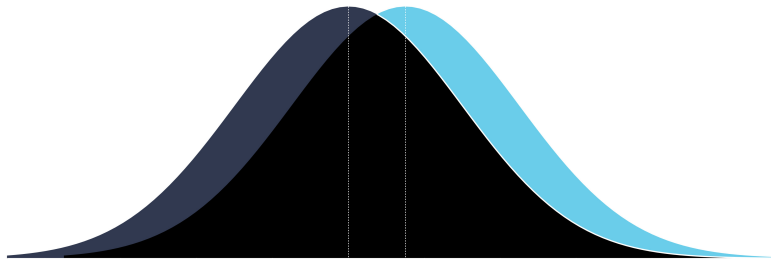


Figure 2.3: A graphical representation of 0.5 Cohen's d (Magnusson, 2023)

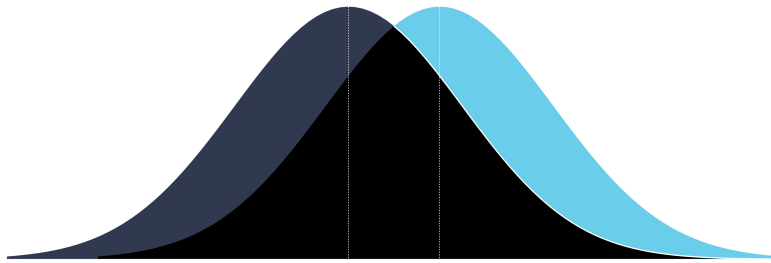


Figure 2.4: A graphical representation of 0.8 Cohen's d (Magnusson, 2023)

While Cohen himself intended these benchmarks as a last resource reference, psychology research has adopted them as standards. Unfortunately, applying these benchmarks without considering the context leads to uninformative conclusions (Funder & Ozer, 2019; Lakens, 2016). Effect sizes require

different interpretations depending on the field of interest, the specific content, and the research methods used, as in a given field a small effect size can be extremely relevant in practice (Altoè et al., 2020).

2.4 Alternative approaches to interpret effect sizes

This clarification has motivated the development of more appropriate benchmarks for effect sizes that are more context specific and informative.

A first category of solution proposed by Funder and colleagues (2020), is comparing findings to well-established results in the literature, e.g. one could compare effect sizes typically considered “small” with classical results obtained in studies on attitude change (Festinger & Carlsmith, 1959), or with to others non-psychology-related phenomena. However, by itself, this approach does not fully resolve the issue. Even if it could be useful to know whether results are comparable to those of another study, if there is no concrete or meaningful sense of the effect in the previous study, we know merely that the effects are similar, but nothing about the meaning they convey. The implications are even less clear when applied across different fields. Therefore, solutions relying solely on conventional or representative benchmarks are unlikely to be sufficient.

To establish more convincing benchmarks though, several methodologies have been proposed. Two of the main approaches are defining a “plausible effect size” (Gelman & Carlin, 2014) and identifying the “smallest effect size of interest” (*SESOI*) (Anvari & Lakens, 2021), both of which will be outlined below.

2.4.1 The plausible effect size

“The plausible effect size refers to what could be approximately the true value of the parameter in the population” (Altoè et al., 2020, p. 5) This concept was introduced to address the challenges of power analysis, or more broadly, design analysis (Gelman & Carlin, 2014), by helping researchers pre-specify an alternative hypothesis, thereby enhancing inferential strength for the reasons previously discussed.

Although the validity of the inferential process is generally strengthened when hypotheses are derived from theory and formalized in statistical terms, precise hypotheses are often not yet feasible in psychology for the reasons previously discussed. Therefore, their formalization can instead be based on literature reviews and/or meta-analyses (Altoè et al., 2020).

The main advantage of plausible effect sizes is to enable a critical interpretation of the results, effectively discriminating plausible from implausible results. With this methodology, we can avoid adjusting effect sizes, either implicitly or explicitly, to justify the value of a given sample size, as it’s often done (Gelman & Carlin, 2014). As stated by Altoè “in general when observed effect size falls outside the pre-specified plausible interval, we can conclude that the observed study is not coherent with our theoretical expectations. On the other hand, we could also consider that our plausible interval may be unrealistic and/or poorly formalized.” (2020, p.9)

Although this methodology is useful to establish criteria that make us cautious about implausible effects, its direct impact on the inferential process is limited, as they do not directly provide with guidelines to distinguish which effect sizes should be deemed meaningful within a given line of research. In exploring further alternative solutions to this issue, the second aforementioned approach will now be explained.

2.4.2 The smallest effect size of interest

In the context of clinical research, scientists are particularly interested in the practical significance of a certain quantity of change. The first relevant ideas on this matter have emerged from the medical field, where quantifying the effect of a given phenomenon is clearly essential. These methods aim to quantify a *Minimal Important Difference* (MID). Subspecifications of this concept include the *Minimally Detectable Difference* (MDD), which emphasizes the ease of detecting a given difference, and the *Clinically Important Difference* (CID), which is based on relevant clinical outcomes such as recurrence or risk of rehospitalization (Norman et al., 2003).

One similar concept has been recently proposed in psychology, the aforementioned *Smallest Effect Size Of Interest* (*SESOI*). The *SESOI* can be defined as “the smallest change that is needed in the outcome measure for people to subjectively notice and report a difference in how they feel” (Anvari & Lakens, 2021, p. 1). This concept translates the notion of an “important” change, defined as one that is both minimal and noticeable, into psychological research, providing a practical criterion for assessing the relevance of effect sizes.

The idea is to assess the importance of a given amount of change based on an external practical criterion. As previously discussed, such a criterion cannot yet rely on solid theoretical reasoning (Riesthuis et al., 2024). However, several methodologies have been proposed, though still being at an exploratory stage. The most prominent among these include cost-benefit analyses, anchor-based methods, and consensus methods (Anvari & Lakens, 2021).

- 1) Cost-benefit analysis evaluates whether the observed improvement justifies the cost of the intervention, usually in comparison with alternatives. However, as noted by Anvari and Lakens themselves, “in basic psychology research, costs and benefit are not easily quantified” (2021, p.2).

- 2) Anchor based methods use a retrospective judgment as a reference to determine whether participants have improved, stayed the same, or worsened over some period of time (Anvari & Lakens, 2021; Lydick & Epstein, 1993; Norman et al., 2003). For example to estimate how much change has been experienced after a treatment, participants complete the measurement of interest before and after the treatment. They are then asked to indicate how much change they notice. Responses may vary depending on the rating scale used. These ratings can be collected through self-reports or clinician’s assessments, and are called clinical “anchor”. Then, the amount of change (effect size) corresponding to each category is measured (Anvari & Lakens, 2021). However, these methods are quiet recent, and not very much established.
- 3) Consensus methods are a recently proposed solution that involve asking experts for their opinion on what could constitute the smallest effect size of interest. Researchers then assess whether there is general agreement on the *SESOI* (Anvari & Lakens, 2021; Riesthuis et al., 2022).

Together, these methods provide preliminary yet important frameworks for defining practically meaningful effect sizes in psychological research.

2.5 A methodology for building a plausible *SESOI*

Based on what has been explained so far, it should be clear how a well-established and pre-specified effect size can strengthen inferential process and enhance the overall credibility of psychological research. In particular, the concept of a plausible effect size provides a solid basis for defining reasonable intervals for the true effect of a phenomenon. Conversely, methods for determining a Smallest Effect Size of Interest (*SESOI*) offer valuable guidance for identifying effects that are

practically meaningful.

In table 2.1 we provide a concise summary of the advantages and disadvantages of both approaches.

Table 2.1: A synthetic comparison of the *SESOI* and the plausible effect size.

	Plausible effect size	SESOI
Pros	• Evidence based	• Practically informative
Cons	• Practically uninformative	• Strongly subjective
Require	• High-quality meta-analysis	• Attentive elicitation procedures

Building on this foundation, we propose a methodology that uses both concepts in a complementary way, leveraging their respective advantages. Specifically, we seek to assess whether the *SESOI* identified is contained within the plausible interval for the effect size or, conversely, whether the effect sizes considered plausible include effects that are practically meaningful, thus ensuring that the estimated true effect is not only statistically valid but also relevant.

The procedure goes as follows:

- 1) Since systematic reviews and meta-analyses can provide guidance on typical effect sizes (Gelman & Carlin, 2014), we can derive a plausible effect size using meta-analytic data. For illustration purpose, we assumed hypothetical meta-analytic results, evaluating efficacy of a clinical intervention, measured with a questionnaire that asses well-being levels. In this scenario, the meta-analysis indicates that the Coehn’s d compared to the control was 0.44, with a 95% confidence interval from 0.27 to 0.61, indicating an improvement in well-being levels. We decide to set our plausible effect size within the confidence interval. A graphical

hypothetical representation is shown in figure 2.5.

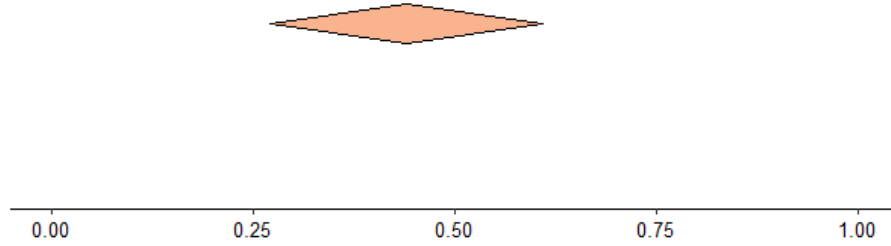


Figure 2.5: A graphical representation of a hypothetical plausible effect size.

- 2) As consensus methods are appropriate for guiding decisions in setting a *SESOI* that reflects practical relevance (Riesthuis et al., 2024), we will use this approach to generate one. We will give particular attention to the matter of consensus methods in the following chapter. For clarity, we generated three different possible *SESOI* scores that may emerge from three different elicitation processes, as graphically represented in figure 2.6.

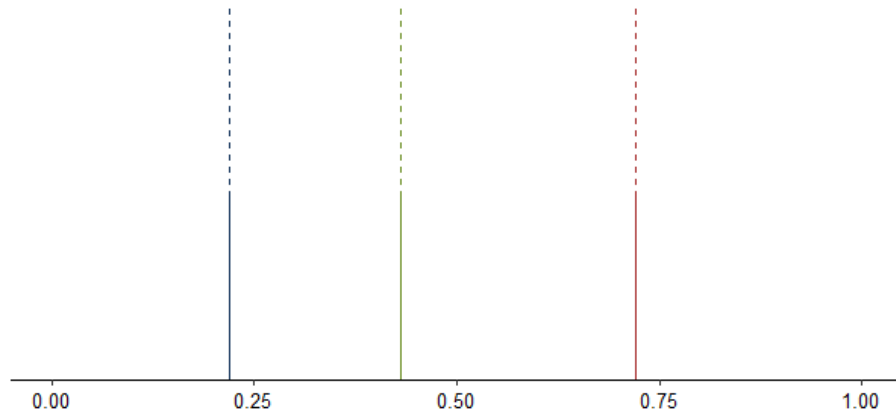


Figure 2.6: A graphical representation of the hypothetical *SESOIs*.

- 3) The two results are then compared, as illustrated in figure 2.7 and discussed. The matter of comparison requires critical and deep thinking, depending heavily on the choices about how both the *SESOI* and the plausible effect size are collected (we will provide a concrete example of such choices and discuss the direct implementation of this framework in chapter 4.) Following the comparison, results are interpreted. In the case of the first and smallest generated *SESOI* score, for example, we might think that since the plausible effect size falls after the line of clinical significance, the vast majority of results in this particular field of research are clinically relevant. With the second and middle-ranged *SESOI* score, we could say that most of the results in this field are somewhat clinically relevant. With the last and largest *SESOI* score, we might argue that our clinical interventions are clinically irrelevant, even though the effect size might sometimes be even medium or large, based on Cohen's standards. Conversely, we could also say that our questionnaire are not suited to effectively identify a clinically relevant change. Obviously, results are never so easily discussed, and deeper reasoning is yet to be done.

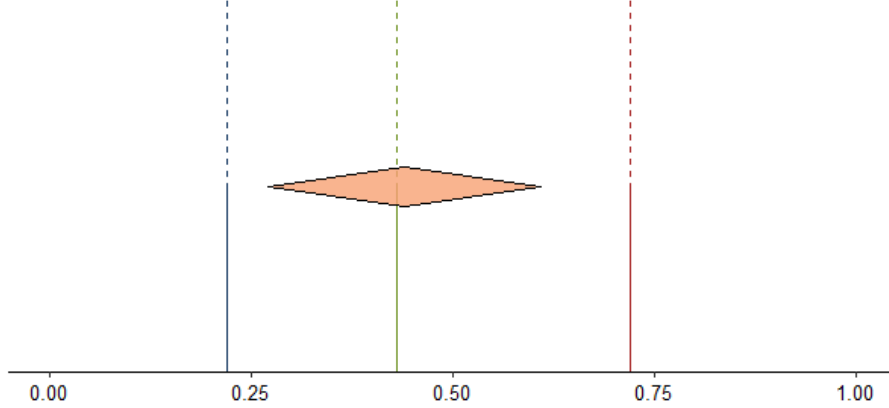


Figure 2.7: The comparison between the plausible effect size and the *SES*Os.

- 4) Finally, a design analysis will be conducted on the obtained result, to illustrate the practical utility of the proposed methodology.

As consensus methods have only recently and infrequently been used to assess a *SESOI*, the following chapter will focus on these approaches, and more specifically, on their more systematic form: expert elicitation.

2.6 Conclusions

In the previous chapter, we discussed how the *NHST* framework often leads researchers to overlook power, sample size, and the practical meaning of results, with p-values frequently misused as indicators of effect strength. Here, we focus on effect sizes as essential tools for addressing these issues, as they connect directly to both study planning and result interpretation.

We defined the effect size, outlined common interpretative pitfalls, and distinguished statistical from practical significance. To enhance interpretation, we presented two complementary approaches: the plausible effect size, based on prior evidence or theory, and the Smallest Effect Size of Interest

(*SESOI*), based on practical relevance. We then proposed a method to integrate both, comparing whether plausible effects are also meaningful.

This chapter shows that careful consideration of effect sizes strengthens inference, supporting the application of the Neyman-Pearson framework, and improves research utility. In the next chapter, we will explore expert elicitation as a method to define the *SESOI* within our framework.

Chapter 3

Expert Elicitation

"Expert opinion and judgment enter into the practice of statistical inference and decision-making in numerous ways. Indeed, there is essentially no aspect of scientific investigation in which judgment is not required.."

— O'Hagan, 2019, p.1

In this chapter we will examine expert elicitation, providing a definition and a brief explanation. We will then outline its main challenges, addressing issues related to expertise, estimate quality, cognitive biases, and aggregation methods. This distinction will allow us to introduce the leading protocols for expert elicitation, explaining how structured approaches can improve accuracy and decision-making. We will then discuss the potential application of these methods for establishing a *SESOI*. Finally, we will introduce a tailored implementation of a remote elicitation.

3.1 Why we need expert elicitation

In the previous chapter, we focused on effect size interpretation as a crucial aspect for addressing sample pre-planning and H_1 specification. Having distinguished between statistical and practical significance, we outlined how the numerical value of the effect size alone does not inform us directly about any theoretically nor practically meaningful implication. To enhance interpretation, we then proposed a method to integrate the plausible effect size and the Smallest Effect Size of Interest (*SESOI*). As explained, the two are complementary, as the plausible effect size informs about the statistical plausibility of our findings while the *SESOI* could give us a precious practical insight.

As we have already outlined, *SESOI* can be obtained based on different methods, i.e. cost-benefit analyses, anchor based methods and consensus methods. Having already discussed the limitations of the first two in the previous chapter, among these methods, consensus methods stand out as particularly promising for our framework.

Firstly, by relying on consensus methods we may ensure the obtained effect size is linked to its practical or clinical relevance. In practice, clinical relevance is often determined by expert clinicians' judgement, where clinicians must take critical decisions. We argue that, despite their inherent limitations, these experts constitute the best available competence we have. In practice, since relying on experts' judgement has been our reliable approach so far, it represents the more solid ground we have on which to establish our *SESOI*. This is especially true as long as we find a way to systematize said opinions, and until more precise and accurate criteria are established. Secondly, consensus methods procedures are well known and extensive literature has already been produced on this subject, making its implementation both feasible and sound.

While at a first glance it may be questionable to rely on experts' judgement, in fields where a value is not known, or decision making is crucial, it is reasonable to rely on experts' opinions

(Authority, 2014b). This is even more reasonable since one’s subjective opinion can be collected in a well established procedure, which is mostly known as Expert Knowledge Elicitation (EKE) (O’Hagan, 2019).

In this chapter we will provide an explanation of such methodologies. We will then introduce the protocol we will be using to collect the *SESOI* in chapter four.

3.2 What is expert elicitation

Expert Knowledge Elicitation (EKE) is a well-established and sound procedure that incorporates expert judgement into formal analyses. The literature on EKE is extensive and covers many fields, including statistics, management sciences, economics, and environmental sciences (O’Hagan, 2019).

Generally speaking, *elicitation* is the process of gathering information needed for a given reason (Authority, 2014a). In EKE, this process refers specifically to obtaining information from one or more experts (Authority, 2014a; O’Hagan, 2019), and in research context, elicitation is usually used to obtain and formalize expert knowledge in the form of probability distributions or estimates (O’Hagan et al., 2006a; O’Hagan, 2019).

In this section, we will address four interconnected issues in EKE: the matter of expertise, the quality of estimates, cognitive biases, and the aggregation problem. First, the definition and selection of expertise determine the validity of the elicitation process, as only well-qualified experts can provide meaningful judgements. Second, the notion of a “good” estimate is examined through different criteria (accuracy, bias, calibration, and reliability). Third, since research has shown that expert opinion is not immune to systematic distortions, we will address the matter of biases. Finally, the aggregation problem is discussed, as expert elicitation usually involves multiple contributors but ultimately requires a single output for decision-making.

3.2.1 The matter of expertise

The definition of “expert” and “expertise” is a key issue in EKE. Various definitions exist, depending on the purpose, methodology, and underlying theories (Authority, 2014a; Burgman, 2015; O’Hagan et al., 2006b), but there is no clear or universally accepted distinction between lay and expert judgment (Burgman, 2015). Even though the debate on this matter is ongoing, when running a study it is essential to clearly specify the type of expertise required and the criteria for expert selection (Authority, 2014a). Linked to this issue, another significant issue revolves around the number of experts involved in the elicitation process. From a theoretical perspective, involving as many experts as possible is advantageous, as it improves the accuracy of estimates for the variables of interest. However, empirical evidence suggests that including too many experts becomes problematic. While the ideal number depends on the field and study, a common recommendation is to involve a minimum of 5 experts, with an upper limit typically between 8 and 15 (Authority, 2014a). Using too few will limit the diversity of perspectives and the informativeness of the results, while involving too many will add little value to the final estimates, while excessively increasing cost and time requirements (Authority, 2014a).

3.2.2 The matter of estimates

It is common practice, and almost common sense, to seek expert opinion when facing uncertainty or lack of knowledge on a particular issue. EKE formalizes this reasonable behavior by asking experts estimates for a given quantity, therefore we shall wonder what defines a “good” estimate. There are several criteria which are used to evaluate expert estimates, but the most typically used are accuracy, bias, calibration, and reliability (Burgman, 2015) :

- *Accuracy* is how close the expert’s quantitative estimate is to the true value.

- *Bias* measures a consistent deviation from the true value, in one direction.
- *Calibration* is a measure of how often the expert’s esteemed intervals contain the true value.
- *Reliability* addresses the characteristic of the expert, as it reflects the degree to which their estimates are repeatable and stable over time.

Not all of these factors are equally important, and the balance depends on the specific problem at hand, as in some circumstances some criteria could be more relevant than others (Burgman, 2015). Although there is an open discussion about how accurate and reliable expert judgements are compared to non-expert opinions, evidence shows that experts generally provide better estimates within their area of expertise. For instance, reliability in EKE is often criticized, as judgements are not very stable over time, which reflects some of the compromises one must accept when using expert elicitation (Burgman, 2015). Nevertheless, we rely on expert advice when making decisions in situations where we lack sufficient information, and when it is our best or only source available (Burgman, 2015). This is particularly relevant in fields like clinical psychology, where objective benchmarks are often unavailable and researchers often rely heavily on unstructured clinical judgment. In these situations, expert intuition frequently serves as the primary basis for crucial decisions. Formalizing this necessary but subjective judgment through EKE provides a systematic way to reduce arbitrary and naïve conclusions, making it a crucial methodology for robust clinical research.

3.2.3 The aggregation problem

As we already mentioned, in most cases, judgements are collected from multiple experts, as it increases the quality of the estimated value, but we typically require a single value as the final output for decision-making. This challenge is referred to as the aggregation problem (O’Hagan, 2019).

To address the aggregation problem, two main strategies are commonly employed, the “mathematical aggregation” and the “behavioural aggregation” (O’Hagan, 2019). The first solution is “mathematical aggregation”, which involves eliciting individual judgements from each expert and fitting a probability distribution to each of them. These distributions are then combined into a single aggregated distribution by a mathematical formula. Such formula is known as a pooling rule. Several pooling rules exist, and choosing the most appropriate one is the most significant choice when implementing this approach (O’Hagan, 2019). The second solution is the behavioural aggregation, which relies on interaction among experts. Firstly, experts discuss their views and reach a consensus. Then, a distribution is fitted to the group’s collective judgment. This method, however, is so sensitive to interpersonal dynamics that may introduce bias. One example could be that dominant personalities can influence group discussions, even unintentionally, compromising the final result (O’Hagan, 2019). One of the most significant risks in behavioural aggregation was first noticed by Janis (1972), as the phenomenon of “groupthink”. In this phenomenon, the desire for consensus overrides critical evaluation, leading to sub-optimal decisions.

3.2.4 The matter of biases

Such a reliance on experts’ judgement makes it crucial to understand and consider potential biases that may influence expert decisions. Unstructured or naïve questioning can introduce unwanted cognitive biases, which can affect experts’ judgements (O’Hagan et al., 2006b). Psychological research has identified ways in which superficial or unstandardised questioning can induce cognitive biases in expert judgment (O’Hagan, 2019). Below, we outline some of the most significant ones to provide an overview of the attention this issue has received in the literature:

- **Anchoring:** When asked to provide a numerical estimate, individuals tend to rely heavily

on the initial value they consider (the “anchor”), and successive adjustments remain biased toward that point (O’Hagan, 2019).

- **Availability:** Events that are more easily recalled or more memorable are often judged more probable, which can lead to an overestimation of dramatic or recent occurrences (O’Hagan, 2019).
- **Range-Frequency:** When the possible values of an uncertain quantity are divided into categories, experts tend to distribute probabilities quite evenly across those categories, regardless of the reasonable likelihoods (O’Hagan, 2019).
- **Overconfidence:** Experts often display overconfidence, which may result in confidence intervals that don’t include the true value, therefore leading to low calibration values. This could be due to social pressures to demonstrate expertise, or due to the nature of elicited questions, which often involve less routine or familiar quantities. Sadly, habitual heuristics that serve experts well in everyday tasks may not be as effective in these elicitation contexts (O’Hagan, 2019).

To reduce the influence of these biases, questions in EKE are carefully designed and preplanned to minimise their impact (Authority, 2014a). To make expert knowledge as objectively as possible, elicitation must be carefully structured. This has led to the development of formal protocols, that are essential to enhance reproducibility and transparency while reducing the influence of biases. We will discuss said protocols in the following section.

3.3 The leading protocols

The decision to adopt either mathematical or behavioural aggregation methods has a direct impact on the conclusion we derive from expert judgements. Moreover, it gives even more importance to

structuring proper protocols, which aim to reduce cognitive biases and enhance transparency and reproducibility throughout the process. These protocols differ greatly in their level of formalization, their approach to uncertainty, how they structure expert interaction and the aggregation process itself. Some of these have emerged as the leading ones in the literature. However, it is not possible to determine definitively which protocol is “the best” in terms of accurately capturing experts’ knowledge and beliefs in the form of probability distributions. Even in the rare scenario where the true values of the quantity of interest are eventually revealed, drawing reliable comparisons would require a substantial number of experts, randomly assigned to different protocols, and a large enough number of elicitation tasks across multiple scenarios for appropriate replication and generalization. Nonetheless, the field of expert knowledge elicitation remains a fertile ground for innovative and systematic research efforts to advance its development (O’Hagan, 2019).

EFSA recommends the use of three main elicitation protocols (Authority, 2014a).

- 1) The **Cooke protocol** is based on mathematical aggregation(Cooke, 1991). It follows a structured approach known as the classical model. In this method, before providing estimates for the quantity of interest, experts first provide independent judgments on a set of *seed variables*. These seed variables are ulterior quantities to be asked, related to the nature of the target variables, but whose true value is known. The expert’s accuracy in judging the seed variables serves as an indicator of the quality of their estimates for the unknown quantities, as this protocol assigns weights to experts based on their performance on the seed questions, and with these weights individual judgments are combined into a single aggregated distribution. The more precise one’s seed estimate, the heavier its weight in the mathematical aggregation (O’Hagan, 2019).
- 2) The **Sheffield protocol** utilizes behavioural aggregation (Oakley & O’Hagan, 2010). It

involves two distinct rounds of expert elicitation. During the first round, experts provide their individual assessments privately to a facilitator. These judgments are then shared and discussed collectively, with the goal of understanding the reasons behind the differences in opinion. Following this discussion, in a second round the group works together to reach a consensus judgment. This consensus is intended to represent the viewpoint of a rational, impartial observer, who should be able to compare objectively the different opinions. The protocol requires a skilled facilitator to guide the process, to ensure to minimize the disturbance of group dynamics and mitigate the biases of the discussions (O’Hagan, 2019).

- 3) The **classic Delphi method** (Linstone et al., 1975; Rowe & Wright, 1999), finally, is mostly used to elicit judgments of uncertainty rather than simple point estimates. This approach combines elements of both mathematical and behavioral aggregation, in a particularly flexible and feasible procedure (Brady, 2015). Experts provide their judgments over two or more rounds, with feedback between rounds, to finally generate a collective response. Since anonymity of individual experts is crucial in this approach, there is limited interaction between experts, allowing some sharing of knowledge, while minimizing social biases risk. After the final round, a pooling rule is applied to combine the experts’ distributions into a single aggregated judgment (O’Hagan, 2019).

Even though these methods typically require a significant investment of time, effort, and financial resources, more agile and feasible versions have been developed. These include methods like the IDEA protocol (Hemming et al., 2018), which streamlines expert elicitation through two estimation rounds with intermediate discussion, enabling flexible remote implementation while maintaining methodological rigor.

3.4 Why eliciting a *SESOI*

As already discussed, elicitation protocols sure have many limitation and issues that must be taken into account. Nonetheless, where little information is available and decisions must be made, relying on the most accurate possible judgement we can elicitate from expert could help improve both research and practice within a given field. The fact that their application has extended to fields like clinical psychology further underscores their utility and validity (Jorm, 2015; Sforzini et al., 2022; Wu et al., 2022)

Guidelines often recommend eliciting values rather than judgments (Hemming et al., 2018). But precisely because of the absence of clear guidance for clinicians, it's arguably even more crucial to establish a shared threshold through a transparent, structured, and deliberative process that minimizes individual biases. The very lack of objectivity of the construct at hand does not weaken the case for expert elicitation, but strengthens it further.

Moreover, we believe that expert elicitation could prove particularly valuable in the interpretation of effect sizes in psychological research, especially given the absence of clear criteria from either theory or other sources, as we outlined in chapter two and three.

3.5 A tailored procedure for eliciting *SESOI*

Since formal methods for a proper elicitation, such as the IDEA protocol, were not feasible within the constraints of this study regarding time and resources, a pragmatic, direct-to-expert approach was developed. This tailored procedure was designed to quickly capture the collective clinical intuition of practicing therapists, translating their experiential knowledge into a quantifiable benchmark for a minimal, yet meaningful, improvement.

In the following chapter we will detail and implement this tailored elicitation procedure to address the central aim of this thesis: to advance effect size interpretation, within the Neyman-Pearson inferential framework, by providing a principled method for specifying a meaningful alternative hypothesis, and inform decision-making when theoretical guidance is lacking.

Chapter 4

Methodology

In this chapter, we will apply the aforementioned methodology to integrate the plausible effect size and the smallest effect size of interest (*SESOI*), addressing its challenges and limitations. We will then conduct a power analysis for a preregistered study, informed by the outcomes of the previous steps, as a case study to demonstrate the practical utility and implications of this integrated approach for further research.

4.1 Introduction

In the previous chapters, we highlighted the conceptual difference between statistical and practical significance, introducing respectively the concepts of the plausible effect size and smallest effect size of interest (*SESOI*). To bridge this gap, we presented a novel methodology to integrate these two approaches, thereby enhancing the interpretation of research findings and theory building in psychology, with a specific focus on clinical psychology.

To emphasize the importance of rigorous and transparent research, we recommend this method-

ology as a preliminary step for a preregistered study or a registered report. Preregistration is a tool for enhancing research transparency by declaring the research hypotheses, methodology, and analysis plan before data collection (Center for Open Science, 2024; Ummul-Kiram et al., 2021). This practice strengthens inferential validity and leads researchers to conceptualize and critically reason through all study phases prior to its implementation. This a-priori reasoning aligns conceptually with the core aim of this thesis: to help pre-specify an effect size that is both empirically grounded and clinically relevant.

For completeness, we note that while preregistration mitigates issues like HARKing (Wagenmakers et al., 2012), its vulnerability to selective reporting might make more rigorous solutions preferable, such as the registered report format (Nosek & Lakens, 2014; Ummul-Kiram et al., 2021), where peer review occurs before data collection, incentivizing researchers to prioritize methodological standards over specific results. However, this approach demands substantial upfront investment and is not always feasible (Ummul-Kiram et al., 2021).

Although highly desirable, formal preregistration is not strictly required to implement this methodology. A reasoned and pre-specified effect size always strengthens inference, provides a better theoretical foundation, and enables appropriate study power.

In this chapter we will move from the theoretical framework to its practical application, as follows. First, we will derive a plausible effect size from the existing literature and a *SESOI* through a process of expert elicitation, as justified in Chapter 3. We will then outline the methodological and conceptual challenges of comparing these two distinct measures, along with our proposed solution. Subsequently, to ground this methodology in a practical context, we will conduct a power analysis for a preregistered study, informed by the outcomes of these initial steps. Finally, the results from this applied methodology will be presented and discussed.

Before detailing these steps we first introduce the primary measure used in our application.

The Kessler-10 scale

To ground our methodology in a concrete example, we will apply it to the Kessler-10 (K-10) (Kessler et al., 2002), a screening tool for psychological distress. Since the SESOI is inherently measure-specific, what constitutes a minimal important change depends entirely on the instrument used. Therefore it is essential to define and justify the chosen measure.

As psychological distress is an umbrella term encompassing various constructs, such as stress, anxiety, and depression, the K-10 is specifically designed to capture symptoms primarily related to anxiety and depressive disorders. This 10-item self-report questionnaire is widely used in international epidemiological research to identify individuals with likely mental health disorders. Its brevity, strong psychometric properties (Lace et al., 2019; Wojujutari & Idemudia, 2024), and free availability make it a practical tool for both research and primary care settings. Respondents indicate how frequently they experienced each symptom over the past 30 days on a 5-point Likert scale (from 1 = “None of the time” to 5 = “All of the time”). Total scores range from 10 to 50, with higher scores reflecting greater severity of psychological distress.

4.2 Methods and materials

In the following subsections we will detail the specific procedures applied for each phase, in the following order:

- Meta-analysis: we will select a meta-analysis from literature to determine the plausible effect size.

- Elicitation: we will describe a tailored elicitation procedure applied to derive the smallest effect size of interest (*SESOI*).
- Comparison: we will address the challenges linked to the comparison of these distinct measures and present one possible solution.

4.2.1 Meta-analysis

We selected a meta-analysis by Madrid et al. (2025) to build up an illustrative case. This recent work synthesizes evidence on digital mental health interventions for university students with mental health difficulties, investigating outcomes across both anxiety and depression domains.

For the present study, we extracted the effect size for depression outcomes from this meta-analysis to serve as our plausible effect size. This choice aligns with the expert elicitation procedure, which was about the score variation of the K-10 using a clinical vignette featuring a patient with depressive symptoms.

The effect size we extracted was $d = 0.55$

4.2.2 Elicitation procedure

We will now outline how the survey was constructed, how experts were recruited and data collected.

Since time and resource constraints prevented us from conducting a manualized elicitation, we conducted a tailored elicitation, that could best meet our resources limitation. The elicitation consisted of a survey via e-mail, with a one-shot question to answer. The choice of this non-manualized methodology is considered appropriate within the context of a master's thesis, intended for illustrative purposes only, and serves as a substitute for a proper elicitation process. The complete elicitation email is reproduced in its original language in Appendix A.

Survey structure

The survey was centered around a detailed clinical vignette depicting a university student with emerging depressive symptoms. The vignette was designed to illustrate a scenario where, after an initial treatment, the patient exhibits a minimal, yet meaningful, functional improvement. Neither length of the treatment nor type of treatment was detailed. This choice may be discussed, but it's our opinion that given the very definition of *SESOI* (the smallest effect size of interest) we are not interested in the time needed to obtain a relevant improvement, nor we are interested in how the improvement is obtained (whether it may be due to a certain treatment, or even spontaneous). Conversely, we are only interested in the score difference, given the relevant change.

Following the vignette, the survey presents the pivotal question: "Given such a relevant change, what difference in the total K-10 score would you expect?" (the original and complete version can be found in the appendix).

For illustrative purpose, we provide a translation of the complete survey above. Should this protocol be used in other contexts, we recommend appropriate adaptation, as the original was administered in Italian.

ELICITATION EMAIL SCRIPT

Subject: Participation in Master's Thesis Study

Body:

Dear [Dr. X],

Thank you for agreeing to take part in this data collection. The results will be used for a

Master's thesis project on methodological research in clinical psychology.

Completion time: **<5 min**

Our goal is to determine what change in the Kessler-10 scale score indicates a minimal clinically significant improvement, i.e., a first real progress towards better functionality for the patient.

We understand that quantifying your opinion might be difficult, but at this stage of the study, it is essential to rely on the clinical experience of professionals like you.

Please consider the following clinical case as an example.

The patient is a university student, with a family history of depression and a history of bullying during his developmental years. He has recently moved to a new city to begin his studies but is encountering academic difficulties. Noticing that the patient is beginning to show a depressive symptomatology, you administer the Kessler-10.

The patient then undergoes a course of treatment. After the treatment, you detect a minimal improvement in the patient's functioning.

Consequently, you re-administer the Kessler-10.

What total score difference do you expect between the first and second administration?

I expect a difference of [x] points.

If you have any doubts, please still indicate the value you consider most plausible. At this stage of the study, your personal "clinical intuition" is the data we are interested in capturing.

Your responses will be used exclusively for research purposes and will be shared only in aggregate form.

We kindly ask you to reply to this email within one week (Wednesday, October 29).

You will receive a reminder on the closing day.

Thank you for your valuable contribution.

Best regards,

Emanuele Bollini - Master's Degree Student

Gianmarco Altoè - Supervising Advisor

The methodology was inspired by the final two stages of the IDEA protocol (notably, “estimate” and “aggregate” phases), which we detailed in the previous chapter (Hemming et al., 2018). In line with the “estimate” phase, experts were asked to provide their best-guess estimates independently and privately, avoiding group dynamics like groupthink, within the one-week timeframe suggested by the protocol. Following the steps of the “aggregate” phase, the final IDEA phase, the collected responses were mathematically aggregated to derive a group estimate.

Pilot test

To refine the clarity and feasibility of the elicitation instrument, a pilot test was conducted. A sample of 6 master's students in clinical psychology was recruited. After being provided with a description of the K-10 scale, they were sent the survey. Their feedback and responses were used to identify ambiguities and ensure clarity.

Ethical

Since no sensitive data were collected, ethical committee approval was deemed unnecessary. Participants were assured that their individual responses would be presented exclusively in aggregated form.

Participants

A panel of 8 experts in clinical psychology were recruited. Experts were defined as licensed psychotherapists with over 2 years of clinical experience, who administer or used to administer the K-10 in their clinical practice. Recruitment occurred by personal invitations, based on a combination of convenience and availability

Data collection

Upon agreeing to participate in the data collection, participants were sent the questionnaire directly via email, without a preliminary explanatory phase. The survey was distributed via email on Wednesday, October 22 at 09:15 to six of the participants, and a second survey was distributed via email on Monday, October 27 at 9:15 to the remaining two (whose email address took longer to collect.) A reminder was sent on the scheduled deadline day (October 29) at 09:15 for the first survey and on the scheduled deadline day (November 3) for the second round.

Exclusions and interpretation

Out of 8 experts, 3 responses were excluded as missing data for different reasons. One expert requested additional information, deeming the protocol invalid without it. Despite our clarification, the expert failed to provide a quantitative estimate. One expert raised concerns about the clinical relevance of the parameter proposed in the answer, casting doubt on the reliability of this response. One expert, during the response phase, revealed information that confirmed to have been erroneously included in the panel, not meeting the pre-defined expertise criteria.

Out of five accepted responses, two responses have been discussed. One response offered two different estimate under different baseline scenario, proposing 5 under moderate distress and six-to-

eight under severe distress condition. Since our clinical vignette did not depicted a severe distress condition, and since its coherent with the definition of a minimal yet meaningful improvement, we took the lower response. One response offered six-to-seven as value. Relying on the same reasoning as the previous case, we considered the lower value.

Following these methodological decisions, the resulting five values were “5, 5, 5, 5, 6”.

Aggregation

Since we are dealing with constructs that cannot be verified, we could not operate techniques such as performance weighting, as described in the previous chapter, which require questions with known answers to calibrate expert weights. This constraint made the use of equal-weight aggregation inevitable, a methodological compromise widely accepted in similar contexts, despite carrying the risk of diluting the contribution of the more accurate experts. Although the choice of aggregation method in such methodologies is often subject to discussion, in this case the mode, median, and mean were all equal to 5. Therefore, we established the value of 5 as the reasonable definitive reference value.

4.2.3 Comparing *SESOI* and Plausible effect

As outlined in chapter 2, our primary objective is to juxtapose the *SESOI* with the plausible effect size, as this will enable us to conduct a power analysis informed by the comparison of the two indices. A key methodological challenge lies in the inherent difference between the plausible and the *SESOI*, as these two indices are derived from different study designs. In this section we will illustrate the differences between the two and under which assumption we are able to do the comparison.

Metanalysis

The plausible effect size from the meta-analysis (Madrid-Cagigal et al., 2025) is a between-groups Cohen's d . This metric is derived from studies that use a randomized controlled trial (RCT) design, where participants are randomly assigned into separate groups, i.e. a treatment group and a control group. Participants in the control group are given a placebo treatment, while participants in the treatment group are given the actual treatment. After the intervention, the effect size is calculated as follows:

$$d_{between} = \frac{M_{treat} - M_{control}}{SD_{pooled}}$$

where (M_{treat}) is the average score of the treatment group after the intervention, $M_{control}$ is the average score of the control group, and (SD_{pooled}) is a weighted average of the standard deviations from the two independent groups. In other words, the ($d_{between}$) expresses the post-treatment difference of treatment groups relative to control group, in terms of standard deviations.

SESOI

In contrast, the expert-elicited *SESOI* represents a within-subject pre-post change, where a single group of participants is measured both before and after an intervention. The effect size is calculated as follows:

$$d_{within} = \frac{\bar{X}_{post} - \bar{X}_{pre}}{SD_{diff}}$$

where \bar{X}_{post} and \bar{X}_{pre} are the mean scores of a single group on a measure (e.g., K-10) after (post) and before (pre) an intervention, and the standard deviation of the difference scores (SD_{diff}) is

given by:

$$SD_{diff} = \sqrt{SD_{pre}^2 + SD_{post}^2 - 2r \cdot SD_{pre} \cdot SD_{post}}$$

(Lakens, 2013; Morris & DeShon, 2002), being SD_{pre} and SD_{post} the standard deviation of the scores for the single group before (pre) and after (post) the intervention. In essence, it expresses the average pre-post change within a single group relative to the variation of the group itself.

Assuming equal variances at pre and post, $SD_{pre} = SD_{post} = SD$, the expression simplifies to

$$SD_{diff} = SD\sqrt{2(1-r)}.$$

Where r is the correlation coefficient between the pre-test and post-test scores for the same individuals. Hence, for an expected mean change of $\Delta = 5$ points:

$$d_{within} = \frac{\Delta}{SD\sqrt{2(1-r)}}$$

Empirical studies on the K-10 in university/community samples typically report SD = 7–9 (Andrews & Slade, 2001; Sunderland et al., 2011). Therefore, we will consider two reasonable scenarios, choosing 8 and 9 values, avoiding 7, to model conservative scenarios:

- *SCENARIO A*: If we assume SD = 8 and $r = 0.50$:

$$SD_{diff} = 8 \cdot \sqrt{2(1-0.50)} = 8, \quad d_{within} = \frac{5}{8} = 0.625.$$

- *SCENARIO B*: If we assume SD = 9 and $r = 0.40$:

$$SD_{diff} = 9 \cdot \sqrt{2(1 - 0.40)} = 9.86, \quad d_{within} = \frac{5}{9.86} \approx 0.51.$$

Thus, broadening the range to account for plausible variability, the elicited *SESOI* corresponds to a within-subjects effect size of about d within $[0.50, 0.65]$ under realistic assumptions.

Comparison

The expert-elicited *SESOI* is a within-subject measure, while the design we are planning (an RCT) and the evidence from our meta-analysis are based on between-group comparisons, therefore they are not directly comparable. To inform the power analysis for a standard RCT, the within-subject *SESOI* must be converted into its between-groups equivalent.

Thus, referring to the previous formula,

$$d_{between} = \frac{M_{treat} - M_{control}}{SD_{pooled}}$$

If we make the following assumptions:

- 1) randomization ensures baseline equivalence,
- 2) the control group shows negligible change over time,
- 3) the post-test SD is approximately equal to the baseline SD,

then the expected post-test difference between groups is just the expected within-person improvement in the treatment group, i.e. $\Delta = 5$, but now expressed on the post-test SD scale.

Because

$$SD_{diff} = SD\sqrt{2(1 - r)} \implies SD = \frac{SD_{diff}}{\sqrt{2(1 - r)}}$$

we can substitute this into the between-groups definition and obtain the key conversion (Morris & DeShon, 2002):

$$d_{between} = d_{within} \cdot \sqrt{2(1 - r)}$$

This makes explicit that the difference between within-subjects and between-subjects effect sizes is entirely due to the correlation between pre and post. Following the two scenarios calculated above,

- *SCENARIO A*: If $d_{within} = 0.625$ and $r = 0.50$:

$$d_{between} = 0.625 \times \sqrt{2(1 - 0.50)} = 0.625 \times 1 = 0.625.$$

- *SCENARIO B*: If $d_{within} = 0.51$ and $r = 0.40$:

$$d_{between} = 0.51 \times \sqrt{2(1 - 0.40)} = 0.51 \times 1.095 \approx 0.56.$$

So, under plausible assumptions for K-10 data, widening the interval for plausibility, the between-groups effect size corresponding to the elicited 5-point change is approximately

$$d_{between} \approx 0.55-0.65$$

The obtained result is notably consistent with the meta-analytic estimate $d = 0.55$ reported in the review on digital mental health interventions for university students, thereby informing us that in this scenario the plausible effect size is also clinically relevant, according to our experts elicitation.

Two important considerations follow from this process. First, we must note that this close

alignment is a specific feature of our case, dependent on our particular parameters—the 5-point change defined by experts, the typical K-10 standard deviations, and the pre-post correlations we assumed. Had these elements been different, we would have needed to develop a different rationale for choosing our target effect size. Second, we can be more or less conservative in deciding how close the two values need to be to justify their combined use. This threshold is not fixed and involves a deliberate methodological choice.

4.3 Analysis

We conducted a sensitivity analysis to assess how sample size requirements varied across different effect sizes ($d = 0.50, 0.55, 0.60$) and statistical power levels (70%, 80%, 90%), using the “pwr” package (Champely, 2020) in R (R Core Team, 2023).

The value $d = 0.55$ represents the direct meta-analytic estimate and the lower bound of our converted *SESOI*, while $d = 0.60$ represents the upper bound of our converted *SESOI*. We also included $d = 0.50$ as a more conservative scenario to illustrate the impact of a slightly smaller, yet still reasonably acceptable, effect.

The effect size range was selected based on the plausible values derived in the previous section. We chose an alpha level of 0.05 and a two-independent-groups design for all power calculations.

The R code of the analysis is included in Appendix B.

4.4 Results

We report the sample size requirements from our sensitivity analysis, which systematically examined how participant numbers vary across different statistical scenarios. Table 4.1 displays required

participants per group (total sample size in parentheses) for two-independent-groups t-test with $(\alpha) = .05$, and shows how sample size increases with higher power levels and decrease with larger effect sizes. Each row shows how sample size changes across power levels for a given effect size, while each column shows how sample size varies across effect sizes for a given power level. Values in parentheses represent the total sample size needed across both groups.

Table 4.1: Participants per group across effect sizes and power levels:

Effect Size	Power 70%	Power 80%	Power 90%
$d = 0.50$	51 (102)	64 (128)	86 (172)
$d = 0.55$	42 (84)	53 (106)	71 (142)
$d = 0.60$	36 (72)	45 (90)	60 (120)

4.5 Discussion

Our power analysis is grounded in the matching procedure between the *SESOI* and the plausible effect size, which in this case indicated that it would be reasonable to think that target effect sizes around $d = 0.55$ are both empirically supported and clinically meaningful for sample size determination.

This sensitivity analysis gives valuable insights for study pre-planning by quantifying the trade-offs between statistical robustness and practical feasibility, as illustrated in figure 4.1. This approach can help researchers make informed and sound decisions by pre-specifying effect size expectations and understanding their implications for sample size requirements.

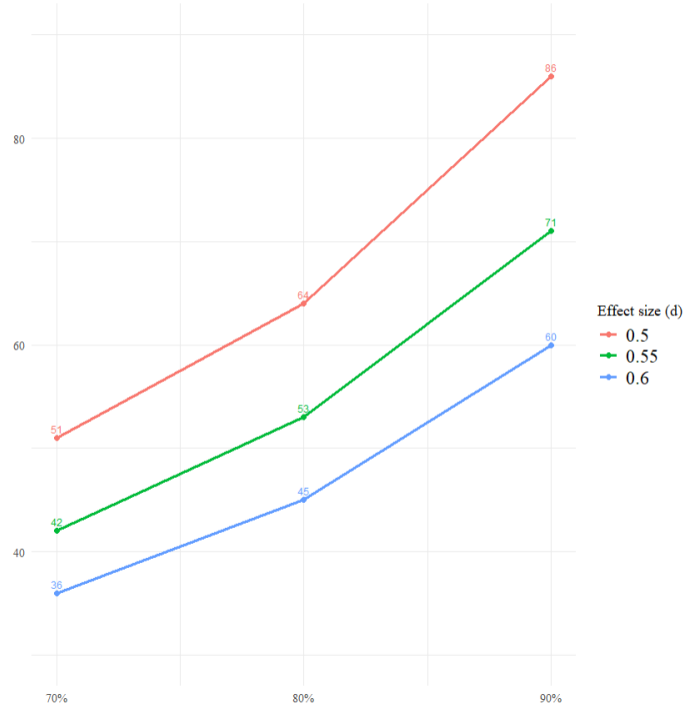


Figure 4.1: Required sample size per group by power and effect size.

As a practical illustration, a researcher adopting the widely accepted power threshold of 80% might pre-specify a conservative target effect size of $d = 0.50$. The application of our methodology provides the justification for this choice, establishing it as both empirically plausible and clinically meaningful. Consequently, the required total sample size would be 128 participants (considering both control and treatment group), a feasible target for many research in psychology.

This approach enables researchers to make informed decisions by pre-specifying effect size expectations and understanding their concrete implications for resource allocation. It is important to acknowledge that this methodological approach has specific limitations, which we will address in detail in the following chapter.

4.5.1 Conclusion

In this chapter we detailed the steps of this novel approach, translating the theoretical methodological framework into an operational procedure. Through a concrete case study, a power analysis for a preregistered study, we have shown how integrating a measure of practical significance (the *SESOI*) with one of statistical significance (the plausible effect size) can directly and meaningfully inform research design.

In the next and final chapter, we will summarize the overarching contribution of this work, discuss its methodological limitations in depth, and propose possible directions for future research to refine.

Chapter 5

Conclusions

Psychology is facing a theoretical and methodological crisis, as discussed in Chapter 1, which is a primary cause for the replicability crisis, undermining the credibility of the discipline. We believe rigorous research practices and transparency are crucial to restore trust in the field.

This thesis therefore aims to strengthen inference moving within the Neyman-Pearson framework. To do this, effect size pre-specification is deemed crucial. One way to encourage effect size pre-specification is to try to develop meaningful interpretation for effect sizes. In order to do so, we proposed and applied a methodology for integrating plausible effect sizes with the smallest effect size of interest (SESOI) (Anvari & Lakens, 2021). The plausible effect size (Altoè et al., 2020) gives the best synthesis literature has achieved about effect sizes on a certain field, while we believe that expert elicitation offers a crucial, clinically grounded interpretation of what constitutes a meaningful change. Combining these two pieces of information greatly helps pre-specifying effect sizes, as it gives a solid line of interpretation for effect sizes, therefore encouraging the adoption of rigorous research standards.

In this concluding section we will synthesize the contributions, acknowledge the limitations, and outline the prospective trajectories of the developed framework.

5.1 General results

The primary contribution of this thesis lies not in a specific empirical finding, but in the development and practical illustration of the procedural framework for enhancing the interpretation and planning of psychological studies. The ultimate aim is to provide a useful methodology for researchers to use and apply in different contexts.

Having addressed the inadequacies in the quality of psychological theories and issues related to the interpretation of effect sizes and the proliferation of under-powered studies, we propose a comprehensive framework designed to guide researchers in conceptualizing and pre-specifying effect sizes, and taking accurate choices in conducting power analysis before the study. This framework systematically combines information from both the existing literature and clinical expertise, creating a more robust foundation for study design.

The applied methodology informed us that a target effect size around $d = 0.55$ should be both clinically relevant and empirically plausible. Building on this information we conducted a power analysis for a range of plausible values around the target effect size.

The sensitivity analysis shows the trade-offs involved in study planning (table 5.1), providing a clear, quantitative basis for resource allocation and study design decisions, moving beyond convention or guesswork.

Table 5.1: Participants per group across effect sizes and power levels:

Effect Size	Power 70%	Power 80%	Power 90%
$d = 0.50$	51 (102)	64 (128)	86 (172)
$d = 0.55$	42 (84)	53 (106)	71 (142)
$d = 0.60$	36 (72)	45 (90)	60 (120)

As discussed in chapter 4, as a practical illustration, a researcher adopting a standard power threshold of 80% and a conservative target effect size of $d = 0.50$, justified by our methodology as both clinically meaningful and empirically plausible, would require a total sample size of 128 participants.

The challenges we encountered largely reflect the pioneering stage of applying SESOI in a clinical context. By documenting the initial exploration of this methodology along with its obstacles, we hope to offer a preliminary step for future research to develop more refined and effective approaches.

5.2 Limitations

While the proposed framework offers a substantial potential, we also acknowledge its limitations, requiring further refinement. The methodological challenges we faced are largely inherent to the novelty of expert elicitation procedures in this context. Rather than strict limitations, these points may be considered constructive pathways for future development, which we will elaborate upon next.

The limitations can be grouped in three main categories:

- **Methodological limitations of the tailored elicitation approach:** These constraints,

primarily concerning the informal nature of our procedure, could be addressed by an established manualized elicitation method or developing a standardized protocol *ad hoc*.

- **Limitations related to the research context:** these limitations derive from a scarcity of necessary data in the existing literature, which constrained the validity of our results. This highlights the urgent need for new primary research to address these gaps and promote open data sharing.
- **Limitations related to statistical assumption:** these constraints were proved necessary in order to permit the comparison. While reasonable, they underscore the need to address the limitations mentioned earlier.

5.2.1 Elicitation related

A key challenge was the inherent circularity in defining the “minimally important change”. To anchor this abstract concept, we used a clinical vignette depicting a patient with emerging depressive symptoms, asking experts to envision a slight but meaningful functional improvement. We then asked the experts to quantify that improvement into the expected Kessler-10 score change. While this provided a concrete reference, the procedure remained intrinsically subjective, relying on individual clinicians’ interpretations of the central construct, potentially introducing more variability. This underscores the importance of eliciting multiple expert judgments to capture diverse clinical perspectives. By eliciting only point estimates rather than uncertainty intervals, this approach does not capture their confidence levels and may reinforce overconfidence bias. This limitation could be mitigated in future research by employing a formal elicitation protocol, such as the IDEA protocol (Hemming et al., 2018), which includes dedicated phases for properly introducing and aligning experts’ understanding of the target constructs. We faced an additional complication when some

experts mistook the SESOI with a measure of treatment efficacy. However, the SESOI defines a clinically meaningful score change on the K-10, regardless of how that improvement is achieved (e.g., treatment, spontaneous remission, or placebo effect). This conceptual confusion was evident when experts inquired about the specific treatment, revealing a focus on the source of improvement rather than the score change that defines it. This underscores a current knowledge gap between methodology and practice. As concepts like SESOI become more integrated into clinical training, this barrier will diminish, facilitating future applications of these techniques.

5.2.2 Research context related

One primary limitation was the heterogeneity of outcome measures across meta-analyses limited the comparability between the meta-analytic effect size and the elicited SESOI. A synthesis based on a single instrument would have been preferable, but such standardization remains rare in the literature. Additionally, regarding the elicitation question, since the literature on treatment efficacy mostly reports effect sizes between-groups using Cohen's d , a directly comparable approach would be to elicit a score change reflecting improvement relative to a control group. We believe, however, that this method may be conceptually problematic for two reasons. Firstly, clinicians lack direct experience with control groups, which could undermine the very expertise elicitation seeks to capture. Furthermore, this framing implicitly assumes that the control group cannot show clinically relevant effects, a significant assumption that must be discussed. While this limits direct comparability, and could be further discussed, our chosen approach prioritizes the ecological validity of the clinical judgment we sought to elicit.

5.2.3 Statistical limitations

As illustrated in previous chapter, the conversion between within-subject and between-groups effect sizes required several statistical assumptions that, while reasonable, represent limitations to our approach. Specifically, we assumed that:

- 1) Randomization ensures baseline equivalence between treatment and control group. This is a fundamental property of randomized designs and ensures that groups are comparable at baseline.
- 2) the control group shows negligible change over time, which is reasonable when using waitlist controls or minimal-attention control conditions that are not expected to produce therapeutic effects.
- 3) the post-test SD is approximately equal to the baseline SD, an assumption commonly made in power analysis and sample size planning when prior information about variance changes is unavailable.

5.3 Future developments

Since our primary goal was to introduce a novel framework, we opted for simplicity as the most effective explanatory approach. Future work, however, could significantly extend this methodology by building more advanced elicitation protocols, collecting more comprehensive data, or even conducting a dedicated meta-analysis tailored for this purpose.

Particularly, building on our experience, we believe that important future developments can be structured around the two main categories of limitations we encountered:

- 1) **Refining the Elicitation Methodology:** future studies should move beyond our tailored

approach by implementing a formal, manualized elicitation protocol, such as the full IDEA framework, or developing a more appropriate variation, to address the issue of the introduction of the SESOI concept, also mitigating the conceptual confusion we observed between a meaningful score change and a specific treatment effect. Following, these approaches could elicit probability distributions or credible intervals instead of single point estimates. As previously outlined, this would capture expert uncertainty, reduce overconfidence bias, and provide a more solid basis for power analysis.

- 2) **Enhancing the Research Context and Integration:** given the constraints of current literature, future studies could conduct a meta-analysis focused on a single, clinically key outcome measure (like the K-10), already adapted for the within metric. This would directly solve the problem of heterogeneous metrics and provide more precise, comparable estimates of the plausible effect size. Lastly, future research could test this framework by designing new studies where SESOI elicitation guides the planning from the very beginning. This would make it possible to collect all the necessary within-subject, avoiding the need of relying on large statistical assumptions and approximations.

Once these current limitations have been adequately addressed, future research could focus on applying this methodology to different questionnaires and psychological constructs.

In conclusion, we hope this framework can encourage researchers to adopt more rigorous practices. Our goal is for it to inspire more critical thinking about effect size and statistical power especially in clinical psychology, where theoretical, methodological and ultimately clinical decisions have a direct impact on people's lives, an ethical imperative that can not be dismissed superficially.

References

- Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J., Gramfort, A., Thirion, B., Varoquaux, G., Accolla, E., et al. (2015). Aarts, AA et al.(2015). Estimating the reproducibility of psychological science. *Science* 349, aac4716–aac4716. DOI. *Modelling Structure and Function of the Human Subcortex*, 52, 299.
- Altoè, G., Bertoldo, G., Zandonella Callegher, C., Toffalini, E., Calcagnì, A., Finos, L., & Pastore, M. (2020). Enhancing statistical inference in psychological research via prospective and retrospective design analysis. *Frontiers in Psychology*, 10, 2893.
- Andrews, G., & Slade, T. (2001). Interpreting scores on the kessler psychological distress scale (K10). *Australian and New Zealand Journal of Public Health*, 25(6), 494–497.
- Anvari, F., Kievit, R., Lakens, D., Pennington, C. R., Przybylski, A. K., Tiokhin, L., Wiernik, B. M., & Orben, A. (2023). Not all effects are indispensable: Psychological science requires verifiable lines of reasoning for whether an effect matters. *Perspectives on Psychological Science*, 18(2), 503–507.
- Anvari, F., & Lakens, D. (2021). Using anchor-based methods to determine the smallest effect size of interest. *Journal of Experimental Social Psychology*, 96, 104159.
- Authority, E. F. S. (2014a). EFSA guidance document for evaluating laboratory and field dissi-

- pation studies to obtain DegT50 values of active substances of plant protection products and transformation products of these active substances in soil. *Efsa Journal*, 12(5), 3662.
- Authority, E. F. S. (2014b). Guidance on expert knowledge elicitation in food and feed safety risk assessment. *EFSA Journal*, 12(6), 3734.
- Banks, G. C., O'Boyle Jr, E. H., Pollack, J. M., White, C. D., Batchelor, J. H., Whelpley, C. E., Abston, K. A., Bennett, A. A., & Adkins, C. L. (2016). Questions about questionable research practices in the field of management: A guest commentary. In *Journal of Management* (1; Vol. 42, pp. 5–20). Sage Publications Sage CA: Los Angeles, CA.
- Borenstein, M., Hedges, L. V., Higgins, J. P., & Rothstein, H. R. (2021). *Introduction to meta-analysis*. John Wiley & Sons.
- Boring, E. G. (1919). Mathematical vs. Scientific significance. *Psychological Bulletin*, 16(10), 335.
- Brady, S. R. (2015). Utilizing and adapting the delphi method for use in qualitative research. *International Journal of Qualitative Methods*, 14(5), 1609406915621381.
- Burgman, M. A. (2015). *Trusting judgements: How to get the best out of experts*. Cambridge University Press.
- Center for Open Science. (2024). *Center for open science*. <https://www.cos.io/>
- Champely, S. (2020). *Pwr: Basic functions for power analysis*. <https://CRAN.R-project.org/package=pwr>
- Cohen, J. (1962). The statistical power of abnormal-social psychological research: A review. *The Journal of Abnormal and Social Psychology*, 65(3), 145.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Routledge. <https://doi.org/10.4324/9780203771587>
- Cohen, J. (1994). The earth is round ($p < .05$). *American Psychologist*, 49(12), 997.

- Collaboration, O. S. (2015). Estimating the reproducibility of psychological science. *Science*, 349(6251), aac4716.
- Cooke, R. (1991). *Experts in uncertainty: Opinion and subjective probability in science*. Oxford university press.
- Derksen, M., & Morawski, J. (2022). Kinds of replication: Examining the meanings of “conceptual replication” and “direct replication.” *Perspectives on Psychological Science*, 17(5), 1490–1505.
- Dienes, Z. (2008). *Understanding psychology as a science: An introduction to scientific and statistical inference*. Bloomsbury Publishing.
- Eronen, M. I., & Bringmann, L. F. (2021). The theory crisis in psychology: How to move forward. *Perspectives on Psychological Science*, 16(4), 779–788.
- Ferguson, C. J. (2016). An effect size primer: A guide for clinicians and researchers. *American Psychological Association*.
- Festinger, L., & Carlsmith, J. M. (1959). Cognitive consequences of forced compliance. *The Journal of Abnormal and Social Psychology*, 58(2), 203.
- Fife, D. A., & Rodgers, J. L. (2022). Understanding the exploratory/confirmatory data analysis continuum: Moving beyond the “replication crisis.” *American Psychologist*, 77(3), 453.
- Fisher, R. (1955). Statistical methods and scientific induction. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 17(1), 69–78.
- Flake, J. K., & Fried, E. I. (2020). Measurement schmeasurement: Questionable measurement practices and how to avoid them. *Advances in Methods and Practices in Psychological Science*, 3(4), 456–465.
- Fried, E. I. (2020). Lack of theory building and testing impedes progress in the factor and network literature. *Psychological Inquiry*, 31(4), 271–288.

- Funder, D. C., & Ozer, D. J. (2019). Evaluating effect size in psychological research: Sense and nonsense. *Advances in Methods and Practices in Psychological Science*, 2(2), 156–168.
- Funder, D. C., & Ozer, D. J. (2020). "Evaluating effect size in psychological research: Sense and nonsense": corrigendum. *Sage Publications*.
- Gelman, A., & Carlin, J. (2014). Beyond power calculations: Assessing type s (sign) and type m (magnitude) errors. *Perspectives on Psychological Science*, 9(6), 641–651.
- Gigerenzer, G. (2004). Mindless statistics. *The Journal of Socio-Economics*, 33(5), 587–606.
- Grissom, R. J., & Kim, J. J. (2012). *Effect sizes for research: Univariate and multivariate applications*. Routledge.
- Hall III, R. P. (2023). Replication and reproducibility and the self-correction of science: What can JID innovations do? *JID Innovations*, 3(3), 100188.
- Head, M. L., Holman, L., Lanfear, R., Kahn, A. T., & Jennions, M. D. (2015). The extent and consequences of p-hacking in science. *PLoS Biology*, 13(3), e1002106.
- Hemming, V., Burgman, M. A., Hanea, A. M., McBride, M. F., & Wintle, B. C. (2018). A practical guide to structured expert elicitation using the IDEA protocol. *Methods in Ecology and Evolution*, 9(1), 169–180.
- Hilgard, J. (2021). Maximal positive controls: A method for estimating the largest plausible effect size. *Journal of Experimental Social Psychology*, 93, 104082.
- Ioannidis, J. P. (2005). Why most published research findings are false. *PLoS Medicine*, 2(8), e124.
- John, L. K., Loewenstein, G., & Prelec, D. (2012). Measuring the prevalence of questionable research practices with incentives for truth telling. *Psychological Science*, 23(5), 524–532.
- Jorm, A. F. (2015). Using the delphi expert consensus method in mental health research. *Australian & New Zealand Journal of Psychiatry*, 49(10), 887–897.

- Kelley, K., & Preacher, K. J. (2012). On effect size. *Psychological Methods*, 17(2), 137.
- Kerr, N. L. (1998). HARKing: Hypothesizing after the results are known. *Personality and Social Psychology Review*, 2(3), 196–217.
- Kessler, R. C., Andrews, G., Colpe, L. J., Hiripi, E., Mroczek, D. K., Normand, S.-L., Walters, E. E., & Zaslavsky, A. M. (2002). Short screening scales to monitor population prevalences and trends in non-specific psychological distress. *Psychological Medicine*, 32(6), 959–976.
- Lace, J. W., Greif, T. R., McGrath, A., Grant, A. F., Merz, Z. C., Teague, C. L., & Handal, P. J. (2019). Investigating the factor structure of the K10 and identifying cutoff scores denoting nonspecific psychological distress and need for treatment. *Mental Health & Prevention*, 13, 100–106.
- Lakens, D. (2013). Calculating and reporting effect sizes to facilitate cumulative science: A practical primer for t-tests and ANOVAs. *Frontiers in Psychology*, 4, 863.
- Lakens, D. (2016). Improving your statistical inferences. -.
- Lakens, D. (2017). Equivalence tests: A practical primer for t tests, correlations, and meta-analyses. *Social Psychological and Personality Science*, 8(4), 355–362.
- Lakens, D. (2019). The value of preregistration for psychological science: A conceptual analysis. , 62(3), 221–230.
- Lakens, D. (2021). The practical alternative to the p value is the correctly used p value. *Perspectives on Psychological Science*, 16(3), 639–648.
- Lin, M., Lucas Jr, H. C., & Shmueli, G. (2013). Research commentary—too big to fail: Large samples and the p-value problem. *Information Systems Research*, 24(4), 906–917.
- Linstone, H. A., Turoff, M., et al. (1975). *The delphi method* (Vol. 1975). Addison-Wesley Reading, MA.

- Lydick, E., & Epstein, R. (1993). Interpretation of quality of life changes. *Quality of Life Research*, 2, 221–226.
- Machery, E. (2020). What is a replication? *Philosophy of Science*, 87(4), 545–567.
- Madrid-Cagigal, A., Kealy, C., Potts, C., Mulvenna, M. D., Byrne, M., Barry, M. M., & Donohoe, G. (2025). Digital mental health interventions for university students with mental health difficulties: A systematic review and meta-analysis. *Early Intervention in Psychiatry*, 19(3), e70017.
- Maier, M., & Lakens, D. (2022). Justify your alpha: A primer on two practical approaches. *Advances in Methods and Practices in Psychological Science*, 5(2), 25152459221080396.
- Malich, L., & Munafò, M. R. (2022). Introduction: Replication of crises-interdisciplinary reflections on the phenomenon of the replication crisis in psychology. *Review of General Psychology*, 26(2), 127–130.
- Morris, S. B., & DeShon, R. P. (2002). Combining effect size estimates in meta-analysis with repeated measures and independent-groups designs. *Psychological Methods*, 7(1), 105.
- Murphy, J., Caldwell, A. R., Mesquida, C., Ladell, A. J., Encarnación-Martínez, A., Tual, A., Denys, A., Cameron, B., Van Hooren, B., Parr, B., et al. (2025). Estimating the replicability of sports and exercise science research. *Sports Medicine*, 1–21.
- Neyman, J. (1957). "inductive behavior" as a basic concept of philosophy of science. *Revue De L'Institut International De Statistique*, 7–22.
- Norman, G. R., Sloan, J. A., & Wywich, K. W. (2003). Interpretation of changes in health-related quality of life: The remarkable universality of half a standard deviation. *Medical Care*, 41(5), 582–592.
- Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G., et al. (2015). Promoting an open research culture.

- Science*, 348(6242), 1422–1425.
- Nosek, B. A., Hardwicke, T. E., Moshontz, H., Allard, A., Corker, K. S., Dreber, A., Fidler, F., Hilgard, J., Kline Struhl, M., Nuijten, M. B., et al. (2022). Replicability, robustness, and reproducibility in psychological science. *Annual Review of Psychology*, 73(1), 719–748.
- Nosek, B. A., & Lakens, D. (2014). Registered reports. In *Social Psychology*. Hogrefe Publishing.
- Nosek, B. A., & Lakens, D. (2016). Registered reports: A method to increase the credibility of published reports. *OSF*.
- O’Hagan, A. (2019). Expert knowledge elicitation: Subjective but scientific. *The American Statistician*, 73(sup1), 69–81.
- O’Hagan, A., Buck, C. E., Daneshkhah, A., Eiser, J. R., Garthwaite, P. H., Jenkinson, D. J., Oakley, J. E., & Rakow, T. (2006a). Uncertain judgements: Eliciting experts’ probabilities. *John Wiley & Sons*.
- O’Hagan, A., Buck, C., Daneshkhah, A., Eiser, J., Garthwaite, P., Jenkinson, D., Oakley, J., & Rakow, T. (2006b). The psychology of judgement under uncertainty. *Uncertain Judgements: Eliciting Experts’ Probabilities*, 33–59.
- Oakley, J. E., & O’Hagan, A. (2010). SHELF: The sheffield elicitation framework (version 2.0). *School of Mathematics and Statistics, University of Sheffield, UK* ([Http://Tonyohagan. Co. Uk/Shelf](http://Tonyohagan.Co.Uk/Shelf)).
- Oberauer, K., & Lewandowsky, S. (2019). Addressing the theory crisis in psychology. *Psychonomic Bulletin & Review*, 26, 1596–1618.
- Pek, J., & Flora, D. B. (2018). Reporting effect sizes in original psychological research: A discussion and tutorial. *Psychological Methods*, 23(2), 208.
- Primbs, M. A., Pennington, C. R., Lakens, D., Silan, M. A. A., Lieck, D. S., Forscher, P. S.,

- Buchanan, E. M., & Westwood, S. J. (2023). Are small effects the indispensable foundation for a cumulative psychological science? A reply to götz et al.(2022). *Perspectives on Psychological Science*, 18(2), 508–512.
- R Core Team. (2023). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. <https://www.R-project.org/>
- Riesthuis, P., Mangiulli, I., Broers, N., & Otgaar, H. (2022). Expert opinions on the smallest effect size of interest in false memory research. *Applied Cognitive Psychology*, 36(1), 203–215.
- Riesthuis, P., Mesquida, C., & Cribbie, R. (2024). *Statistical (non) significance ≠ (un) successful replication: The importance of the smallest effect size of interest*. OSF.
- Rowe, G., & Wright, G. (1999). The delphi technique as a forecasting tool: Issues and analysis. *International Journal of Forecasting*, 15(4), 353–375.
- Scheel, A. M. (2022). Why most psychological research findings are not even wrong. *Infant and Child Development*, 31(1), e2295.
- Scheel, A. M., Schijen, M. R., & Lakens, D. (2021). An excess of positive results: Comparing the standard psychology literature with registered reports. *Advances in Methods and Practices in Psychological Science*, 4(2), 25152459211007467.
- Scheel, A. M., Tiokhin, L., Isager, P. M., & Lakens, D. (2021). Why hypothesis testers should spend less time testing hypotheses. *Perspectives on Psychological Science*, 16(4), 744–755.
- Schmidt, S. (2009). Shall we really do it again? The powerful concept of replication is neglected in the social sciences. *Review of General Psychology*, 13(2), 90–100.
- Sforzini, L., Worrell, C., Kose, M., Anderson, I. M., Aouizerate, B., Arolt, V., Bauer, M., Baune, B. T., Blier, P., Cleare, A. J., et al. (2022). A delphi-method-based consensus guideline for definition of treatment-resistant depression for clinical trials. *Molecular Psychiatry*, 27(3), 1286–

1299.

- Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science*, 22(11), 1359–1366.
- Sunderland, M., Slade, T., Stewart, G., & Andrews, G. (2011). Estimating the prevalence of DSM-IV mental illness in the Australian general population using the Kessler Psychological Distress scale. *Australian & New Zealand Journal of Psychiatry*, 45(10), 880–889.
- Ummul-Kiram, K., Silverstein, P., & Moin, S. (2021). Easing into open science: A guide for graduate students and their advisors. *Collabra: Psychology*, 7(1).
- Van de Schoot, R., Kaplan, D., Denissen, J., Asendorpf, J. B., Neyer, F. J., & Van Aken, M. A. (2014). A gentle introduction to Bayesian analysis: Applications to developmental research. *Child Development*, 85(3), 842–860.
- VandenBos, G. R. (2007). *APA dictionary of psychology*. American Psychological Association.
- Vazire, S., Schiavone, S. R., & Bottesini, J. G. (2022). Credibility beyond replicability: Improving the four validities in psychological science. *Current Directions in Psychological Science*, 31(2), 162–168.
- Wagenmakers, E.-J., Wetzels, R., Borsboom, D., Maas, H. L. van der, & Kievit, R. A. (2012). An agenda for purely confirmatory research. *Perspectives on Psychological Science*, 7(6), 632–638.
- Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>
- Wojujutari, A. K., & Idemudia, E. S. (2024). Consistency as the currency in psychological measures: A reliability generalization meta-analysis of Kessler Psychological Distress Scale (K-10 and K-6). *Depression and Anxiety*, 2024(1), 3801950.

Wu, H., Xu, L., Zheng, Y., Shi, L., Zhai, L., & Xu, F. (2022). Application of the delphi method in the study of depressive disorder. *Frontiers in Psychiatry*, 13, 925610.

Appendix A

The following section presents the tailored elicitation script used in this study. As the elicitation was conducted with Italian experts, the communication is presented in its original Italian.

ELICITATION EMAIL SCRIPT

Oggetto: Partecipazione a studio per tesi magistrale

Corpo: Gentile [dott./dott.ssa X],

Grazie per aver aderito a questa raccolta dati. I risultati saranno impiegati per un lavoro di tesi magistrale sulla ricerca metodologica in psicologia clinica.

Tempo di compilazione: <**5min**

Il nostro obiettivo è valutare quale variazione del punteggio della scala Kessler-10 indica un minimo miglioramento clinico significativo, ossia un primo reale progresso nel paziente verso una migliore funzionalità.

Sappiamo che quantificare la sua opinione potrebbe essere difficile, ma in questa fase dello studio è

fondamentale basarsi sull'esperienza clinica di professionisti come lei.

Consideri il seguente caso clinico a scopo esemplificativo.

Il paziente è uno studente universitario, con una storia familiare di depressione ed episodi di bullismo in età evolutiva. Si è da poco trasferito in una nuova città per iniziare gli studi, ma riscontra difficoltà a livello accademico. Notando che il paziente comincia a manifestare un quadro sintomatologico depressivo gli somministra la Kessler-10.

Il paziente viene poi sottoposto a un trattamento. Dopo il trattamento, lei rileva un miglioramento minimo del funzionamento del paziente.

Di conseguenza risomministra la Kessler-10.

Quale differenza di punteggio complessivo si aspetta tra prima e dopo?

Mi aspetto una differenza di [x] punti.

Qualora avesse dei dubbi, la preghiamo di indicare comunque il valore che ritiene più plausibile.

In questa fase dello studio, la sua personale “intuizione clinica” è il dato che siamo interessati a rilevare.

Le sue risposte saranno usate esclusivamente per scopi di ricerca e verranno condivise solamente in forma aggregata.

Le chiediamo cortesemente di **rispondere a questa mail entro una settimana (Mercoledì 29 Ottobre)**. Riceverà un promemoria il giorno della chiusura.

Grazie per il suo prezioso contributo.

Cordiali saluti,

Emanuele Bollini - Laureando magistrale

Gianmarco Altoè - Relatore supervisore

Appendix B

This appendix contains the full source code for analyses and graphs executed in Chapter 4.

Analyses were conducted with R, version 4.1.1 (2021-08-10). Power analysis was conducted with “pwr” package (Champely, 2020) Graphs were created using package “ggplot2” (Wickham, 2016)

Power calculation and graphs

```
#START SESSION

rm(list = ls())                                #Clear all objects from workspace

#POWER ANALYSIS

##Required packages

library(pwr)                                  #Load the package

##Setting parameters for the simulation
```

```

ds      = c(0.50, 0.55, 0.60)           #Multiple effect sizes, where
                                         #0.55 is Choen's d
                                         #plausible effect size
                                         #0.50 and 0.60 are two
                                         #scenarios effect sizes
                                         #from the reasonable range
                                         #of sesoi within effect
                                         #sizes  converted to between

alpha   = 0.05                          #Significance level

powers  = c(0.70, 0.80, 0.90)           #Setting different power
                                         #scenarios (/70%,80%,90%))

##Power analysis function

calc_n = function(d, power, alpha = 0.05) {
                                         #This function calculates
                                         #the sample size, given:

    out = pwr.t.test(d = d,              #effect size
                     power = power,      #desired powers
                     sig.level = alpha,  #significance levels
                     type = "two.sample", #two independent groups
                     alternative = "two.sided") #two tailed test

    ceiling(out$n)                       #round up to nearest integer
}

```

```

n_per_group = sapply(ds, function(d) {          #For each effect size in ds

  sapply(powers, function(pw) calc_n(d, pw, alpha))#sample for ALL power levels
})

##Show Result

results = data.frame(                           #creates a data frame that enlists:

  d_between    = rep(ds, each = length(powers)),    #Repeats each d for all powers
  power        = rep(powers, times = length(ds)),    #Repeats powers for each d
  n_per_group  = c(n_per_group),                    #Flattens the matrix
  total_sample = c(n_per_group) * 2                  #Total participants (2 groups)
)

print(results)                                  #shows table of results

#ILLUSTRATIVE GRAPHS

##Required Packages

library(ggplot2)                                #Load the package

##Grpahs

#Create empty list to store all combinations

```

```

res_list = list()

#Generate all effect size and power combinations
for (d in ds) {
  for (pw in powers) {

    #Calculate required sample size for current combination
    n_pg <- calc_n(d, pw, alpha)

    #Store results in data frame and add to list
    res_list[[length(res_list) + 1]] <- data.frame(
      d_between = d,          #Current effect size
      power = pw,            #Current power level
      n_per_group = n_pg     #Calculated sample size per group
    )
  }
}

#Combine all list elements into single data frame
df <- do.call(rbind, res_list)

#Creating plot for required sample size per group by power and effect size
ggplot(df, aes(x = power, y = n_per_group,

```

```

    group = factor(d_between),      #Separate line for each effect size
    color = factor(d_between))) +  #Different color for each effect size

#Plot elements

geom_line(linewidth = 1) +          #Connect points with lines
geom_point(size = 2) +              #Connect points with lines
geom_text(aes(label = n_per_group),  #Add sample size labels above points
          vjust = -0.6, size = 3, show.legend = FALSE) +

#Axis formatting

scale_x_continuous(breaks = powers,      #Set x-axis breaks
                  #at power levels
                  labels = paste0(powers*100, "%")) + #Format as percentages

#Y-axis limits - set fixed range for consistent appearance

scale_y_continuous(limits = c(30, 90)) + #Fix max and min range

#Labels and titles

labs(
  title = NULL,      #(Title and labels were set null for clarity
                    #once added to the final document)
  x      = NULL,      #Desired power - Axis
  y      = NULL,      #Sample size per group - Axis

```

```

    color = "Effect size (d)"          #Legend title
  ) +

#Theme settings
theme_minimal(base_size = 12) +
theme(
  text = element_text(family = "serif"),      #Set all text to Times New Roman
  legend.text = element_text(size = 15)
)

#-----

#VISUALIZING PLAUSIBLE EFFECT SIZE AND SESOI MATCHING

#Package required
library(ggplot2)

#CREATE DATA FRAMES FOR PLOTTING
effs = data.frame(          #Creates empty data
  x = c(0, 0.22, 0.43, 0.72),    #SESOI positions
  y = c(1, 2, 3, 4)              #Vertical positions for layout
)

```

```

poly_df = data.frame(                                #Coordinates for diamond-shaped polygon

  x = c(0.27, 0.44, 0.61, 0.44),

  y = c(0.50, 0.56, 0.50, 0.44)

)

#PLAUSIBLE ES GRAPH

graph_p = ggplot(effs, aes(x = x, y = y)) +          #SPACE FOR THE GRAPH

  theme_classic() +

  theme(

    axis.line.y = element_line(color = "white"),

    axis.text.x = element_text(size = 10, color = "black", family= "Times"),

    axis.text.y = element_blank(),

    axis.title.x = element_blank(),

    axis.title.y = element_blank(),

    axis.ticks.y = element_blank()

  ) +

  coord_cartesian(xlim = c(0, 1), ylim = c(0, 1)) +

  # Polygon describing the effect size diistribution

  geom_polygon(

    data = poly_df,

    aes(x = x, y = y, group = 1),

    inherit.aes = FALSE,

```



```

    fill = "#f7a072",

    colour = "black",

    size = 0.4,

    alpha = 0.8

)

print(graph_p)

#SESOI GRAPH

graph_s = ggplot(effs, aes(x = x, y = y)) +

  theme_classic() +

  theme(

    axis.line.y = element_line(color = "white"),

    axis.text.x = element_text(size = 10, color = "black", family= "Times"),

    axis.text.y = element_blank(),

    axis.title.x = element_blank(),

    axis.title.y = element_blank(),

    axis.ticks.y = element_blank()

  ) +

  coord_cartesian(xlim = c(0, 1), ylim = c(0, 1)) +

  #layer interni (linee/segmenti)

  geom_vline(xintercept = 0.22, linetype = 2, color = "#1b3a60") +

```

```

geom_vline(xintercept = 0.43, linetype = 2, color = "#6b8e23") +
geom_vline(xintercept = 0.72, linetype = 2, color = "#a52a2a") +
geom_segment(x= 0.22, y = -1, xend = 0.22, yend= 0.5, colour = "#1b3a60") +
geom_segment(x= 0.43, y = -1, xend = 0.43, yend= 0.5, colour = "#6b8e23") +
geom_segment(x= 0.72, y = -1, xend = 0.72, yend= 0.5, colour = "#a52a2a")

print(graph_s)

#GRAPH MATCHED

graph_m = ggplot(effs, aes(x = x, y = y)) +
  theme_classic() +
  theme(
    axis.line.y = element_line(color = "white"),
    axis.text.x = element_text(size = 10, color = "black", family= "Times"),
    axis.text.y = element_blank(),
    axis.title.x = element_blank(),
    axis.title.y = element_blank(),
    axis.ticks.y = element_blank()
  ) +
  coord_cartesian(xlim = c(0, 1), ylim = c(0, 1)) +

  #segments (same code as above)

```

```

geom_vline(xintercept = 0.22, linetype = 2, color = "#1b3a60") +
geom_vline(xintercept = 0.43, linetype = 2, color = "#6b8e23") +
geom_vline(xintercept = 0.72, linetype = 2, color = "#a52a2a") +
geom_segment(x= 0.22, y = -1, xend = 0.22, yend= 0.5, colour = "#1b3a60") +
geom_segment(x= 0.43, y = -1, xend = 0.43, yend= 0.5, colour = "#6b8e23") +
geom_segment(x= 0.72, y = -1, xend = 0.72, yend= 0.5, colour = "#a52a2a") +

#polygon (same code as above)
geom_polygon(
  data = poly_df,
  aes(x = x, y = y, group = 1),
  inherit.aes = FALSE,
  fill = "#f7a072",
  colour = "black",
  size = 0.6,
  alpha = 0.8
)

print(graph_m)

```