

CES in the Treatment of Anxiety Disorders - Part 2

Statistical Considerations in the Meta-Analysis of Cranial Electrotherapy Stimulation (CES) treatment of Anxiety Disorders

By Daniel L. Kirsch, PhD, DAAPM, FAIS, and Marshall F. Gilula, MD



Daniel L. Kirsch, PhD, DAAPM, FAIS
Department Head

While anecdotal results of CES treatment for anxiety disorders are invariably positive, a rigorous, scientific approach is required for analyzing, collating, and reporting results from the vast body of research done on CES. Due to varying methodologies and measures, the myriad of studies do not lend themselves to a simple consolidation of results. Therefore, a statistical method called 'meta-analysis' is used to combine results in a meaningful way and allow an objective measure of the efficacy of CES.

— Daniel L. Kirsch, PhD

— Marshall F. Gilula, MD



Marshall F. Gilula, MD

Part 2 continues from the March 2007 issue of *Practical Pain Management*. Meta-analysis is a statistical method of combining the results of several studies that address a set of related research hypotheses. Because the results from different studies investigating different independent variables are measured on different scales, the dependent variables in a meta-analysis are some standardized measure of effect size. The usual effect size indicator is either the standardized mean difference or an odds ratio in experiments with outcomes of dichotomous variables (success versus failure).

In this case, a meta-analysis of CES calculates the percent of patients improving versus the percent not improving to yield the treatment effect size r , which is equal to the amount of patient improvement given as percentage.³² In the previous issue, it was reported that results of 500 patients produced an effect size $r=.62$. When the smaller groups of patients with specific types of anxiety related disorders were broken out, the effect size among those suffering from panic disorder was $r=.45$, OCD patients, $r=.68$, bi-polar disorder $r=.71$, PTSD ($r=.55$) ADHD ($r=.62$), and phobias ($r=.49$). The overall mean effect size for the combined smaller groups was $r=.64$. These results can be compared with the accepted standardized ratings of $r=.10$ for small effect, $r=.30$ for medium effect and $r=.50$ for large effect.³³ Thus it can be seen that the overall effect of CES for anxiety disorders is large and that there is a notable effect of duration of use that enhances such outcomes.

Statistical Significance

When any given study is published, the authors analyze the data and report whether or not the treatment utilized in their study had a discernable effect. They may report that the treatment had a significant effect at the .05 or .01, or .001 level of probability. In the first instance, the .05 indicates that if the study were to be repeated 100 times, the changes found might have

occurred by chance alone only 5 times out of 100. Or in the case of .01 or .001 level of probability, the result would be expected to have occurred by chance alone only one time out of 100 or one time out of 1,000, respectively.

This form of data analysis and reporting are the hallmarks of contemporary science. Most health care professionals invest meaning in such reporting and deduce that we can have confidence in such data. We can know that the treatment effect is almost certainly genuine and effective, especially if we see one with a .001 probability utilized, as one can assume that a study yielding a probability of $p<.001$ had a really strong clinical effect.

Such considerations are called statistical significance. However, statistical significance does not always tell us anything regarding the actual improvement or efficacy of the treatment studied. For example, what if the study were designed to discern the effect of painting hospital room walls sunlight yellow for severe pain patients? In this hypothetical study, researchers might measure the patients' feelings of well being on a 100 point scale. Suppose most of the patients began at 3 on the scale, with a scoring range from 1 to 5, indicating very low feelings of well being, and went up to 4 on the 100 point scale after their room was painted. Although the average score increased by one-third, and that change was found to be significant at the $p<.01$ level, we are compelled to ask how important is such a finding to the total well being of pain patients, and by extension, what do such results imply in terms of cost and time impact (supposing that one were to use these results to justify painting the walls of hospital wards yellow)?

So how important is statistical significance? The answer can depend on many things, such as how much treatment effect, or patient improvement the treatment yielded, and the significance figure does not provide this. One could also consider what treatment costs are involved in the process of effecting that change, and whether there are other treatments available that can make the same, or even greater changes at less cost. In our hypothetical scenario such factors might involve the cost of scraping the old paint off, removing the mold, repairing and repainting the walls, and comparing that to other treatments that are available

for the same amount of money that might provide equal or greater benefits for pain patients. Most of these questions are not statistical, but are questions of clinical relevance, cost and benefit, and can involve personal values as well.

There is also a second use of the term "significant" in medical literature. Most pharmaceutical companies state that any improvement of 25% or better is significant. Is a 25% improvement also statistically significant? Not necessarily. For example, if we compare the results of a treated group with the results of a sham treated group, the difference may not be statistically significant in that both groups may improve. Both may have improved 25% or more during the course of the study. Such results may not be reported in the journal article or advertisement for a given product. Instead, a statement such as, "40% of the treated group improved significantly at 25% or above" may be all that is provided. So when such studies are read, one needs to always look for any comparison between treated and sham treated subjects. For example, some follow-up studies by public interest groups have shown that several major antidepressant medications were later found to be no better than placebo treatment.³⁴ Also, seasoned neurological researchers opine that many new anticonvulsant drug studies routinely exclude from the treatment group any patients who show an initial intolerance to the drug, and this percentage of a selected population may routinely fluctuate between 12-25% of the population that is selected for testing.³⁵

Effect Size

So while in the past there has been a focus on significant results in scientific studies, we now understand that the term, "significant" can be used in at least two ways. It can be inferred that "significance" alone is no longer an exclusive hallmark of sufficient information. A clinician needs to know how effective a treatment is in terms of the actual amount of improvement it produces in order to make an informed decision about which intervention to use. Certainly there might be less interest in a highly significant statistical result if the reduction of a given symptom is only 3%, and if that were clearly stated in the results section of a journal article or in an advertisement.

If two different studies report the results of two different types of treatments, and

the results of both were found to be significant at the .05 level of significance, one would clearly be more interested in the one that resulted in a symptom reduction of 80% over the one that resulted in a reduction of 15%. This difference is known as the effect size. In advertisements and much of the scientific literature, this is ordinarily not disclosed. The reader cannot know the effect size from a study unless the published results are carefully evaluated for percent improvement pre- to post-treatment, above and beyond that of the controls.

Another problem is that when a treatment is used in studies of various groups in different parts of the country, or with groups showing slight differences in their diagnostic profile (or with groups studied at different times of the year), studies may all report significant improvement of the patients at the .05 level of confidence, but the effect sizes, when these can be ascertained, may vary considerably across the studies. A physician who wants to know what to expect if a medication or device is used in practice cannot accurately derive this knowledge from this type of reporting, and may thus not be able to reproduce the reported effects in actual patients. The best way to determine the overall effect from diverse and numerous studies is through the use of meta-analysis.

Meta-analysis is a statistical technique in which all the effect sizes found in a group of studies of the same treatment can be summarized into an overall average effect size. The derived mean effect size is what one can expect to see in most treated patients, most of the time. If meta-analysis of studies yields an average effect size of 15%, this will be of less interest to the practicing physician than a meta-analytic finding from another treatment which treats the same problem, but results in an average effect size of 60%.

Simply stated, the r effect size represents the percentage improvement to be expected on a scale of 0 to 100. An r effect size of .15 means that there was an average of only 15% improvement among patients when measured across combined studies, while $r=.75$ means that there was an average of 75% improvement in patients found in the combined studies, etc. In this scale, an r effect size of .10 is small, while r of .30 is moderate, and r of .50 or above is considered to be high.

Many early statistical meta-analyses

were confined to studies that specifically reported the pre- and post-study means and standard deviations. All other studies had to be ignored, no matter how rigorous the scientific protocols. That left out the results of some well designed, well conducted double-blind placebo-controlled studies. Current use of meta-analysis tends to statistically transform whatever statistic the author reports into an effect size statistic and then proceed with subsequent analysis from that data set of the collected studies (See Appendix A for an example).

What one might gain from this discussion is that the effect sizes obtained by meta-analytical procedures of CES studies is very robust and holds up to scrutiny very well given the reasonably large number of studies available to work with. The effect size of CES — as derived from Tables 5 and 6 — was seen to stabilize in the high 50s or low 60s, with the expected effect size in 99 out of 100 times in a future meta-analysis of studies to range from $r=.40$ s to $r=.70$ s. That range is considered to represent a moderate to very strong clinical improvement.

Discussion of CES Meta-Analysis Results

The most pristine analysis of the re-ordered data yielded an effect size of $r=.57$ (as opposed to the un-ordered meta-analysis table that yielded $r=.58$). Analysis of the studies that used only the double-blind method provided $r=.53$.

After removing extraneous measures of anxiety and only analyzing for state anxiety or trait anxiety using the State/Trait Anxiety Inventory, $r=.60$ for state anxiety and $r=.68$ for trait anxiety. These involved a relatively small number of studies. When results were corrected for the number of subjects in each study, the r for state anxiety fell back to a more typical $r=.59$, while trait anxiety fell back to a more typical $r=.60$.

What this means is that while Klawansky at Harvard found an average effect size of $r=.53$ in earlier meta-analysis of eight CES studies, and O'Connor in Tulsa found an r effect size of $r=.51$ in the eight CES studies she chose, similar effect size results were obtained when more than five times that number of studies were meta-analyzed. If an additional 400 CES studies of anxiety were to be analyzed 50 years from now, the likelihood is almost certain that the effect size would still be within the

**TABLE 5.
PUBLISHED CES
ANXIETY
STUDIES**

| AUTHOR | NUMBER OF PATIENTS | | | STATISTIC REPORTED* | RESULT | Zr SCORE |
|------------------------------------|--------------------|----------|-------|---|--------|----------|
| | CES | CONTROLS | TOTAL | | | |
| Bianco, 1994 ²⁴ | 29 | 18 | 47 | % Improvement, Beck AI | .77 | 1.02 |
| | | | | % Improvement, Hamilton AS | .71 | .887 |
| Feighner, 1973 ³⁶ | 23 | 23 | 23 | % Improvement | .31 | .321 |
| Flembaum, 1974 ³⁷ | 28 | Historic | 25 | % Pts Much, or Very Much Improved | .51 | |
| Frankel, 1973 ³⁸ | 17 | 17 | 17 | % Improvement | .08 | .080 |
| Gibson, 1983 ⁵ | 16 | 16 | 32 | % Improvement, EMG | .35 | .365 |
| | | | | % Improvement, STAI | .43 | .460 |
| Gomez, 1974 ²³ | 14 | 14 | 28 | % Improvement, TMAS | .35 | .365 |
| Hearst, 1974 ³⁹ | 14 | 14 | 28 | % Pts. Asymptomatic | .71 | |
| Heffernan, 1995 ⁶ | 10 | 10 | 20 | t-score, EMG | 2.35 | .717 |
| | | | | t-score, Heart Rate | 2.55 | .784 |
| | | | | t-score, finger temperature | 2.62 | .717 |
| | | | | t-score, capacitance | 2.14 | .662 |
| Heffernan, 1996 ¹⁶ | 10 | 20 | 30 | % Increase in EEG Correlation Dimension | .54 | .604 |
| Jemelka, 1975 ⁴⁰ | 14 | 14 | 28 | P = < .05 Improvement, Hamilton AS | .51 | .563 |
| Kirsch, 2002 ²⁸ | 298 | | 298 | % Improvement | .83 | 1.188 |
| Krupitsky, 1991 ⁴¹ | 10 | 10 | 20 | % Improvement, State Anxiety | .41 | .436 |
| | | | | % Improvement, Trait Anxiety | .73 | .929 |
| | | | | % Improvement, TMAS | .47 | .510 |
| Levitt, 1975 ⁴² | 5 | 6 | 11 | % Improvement, TMAS | .80 | 1.099 |
| McKenzie, 1976 ⁴³ | 8 | 4 | 12 | % Improvement, Skin Potential | .48 | .523 |
| Magora, 1967 ⁴⁴ | 20 | | 20 | % Improvement | .75 | |
| Matteson, 1986 ¹⁹ | 32 | 22 | 54 | t-score, State Anxiety | 4.63 | .640 |
| | | | | t-score, Trait Anxiety | 3.37 | .523 |
| | | | | t-score, POMS Anxiety | 5.43 | .701 |
| May, 1993 ⁴⁵ | 14 | | 14 | % Improvement, MAACL | .75 | .973 |
| Moore, 1975 ⁴⁶ | 17 | 17 | 17 | % Improvement, Psychiatrist Ratings | .35 | .365 |
| Overcash, 1999 ⁴⁷ | 182 | | 182 | % Change, EMG | .72 | .908 |
| | | | | % Change, Electrodermal Response | .48 | .523 |
| | | | | % Change, Temperature | .13 | .131 |
| | | | | % Change, Self Rating Scale | .76 | .996 |
| Overcash, 1989 ⁴⁸ | 16 | 16 | 32 | % Change, EMG | .92 | 1.589 |
| | | | | % Change, 16PF, Planful Scale | .80 | 1.099 |
| Passini, 1976 ⁴⁹ | 30 | 30 | 60 | % Improvement, MACL | .28 | .288 |
| | | | | % Improvement, State Anx. | .30 | .310 |
| | | | | % Improvement, Trait Anx. | .10 | .100 |
| Patterson, 1984 ⁵⁰ | 186 | | 186 | % Improvement, Anxiety | .75 | .973 |
| Philip, 1991 ⁵¹ | 10 | 11 | 21 | P = < .05 Improvement | .60 | .693 |
| Rosenthal, 1972 ²⁰ | 11 | 11 | 22 | % Improvement | .67 | .811 |
| Rosenthal, 1970 ⁵² | 9 | | 9 | % Improvement | .47 | .510 |
| Rosenthal, 1970a ⁵³ | 12 | | 12 | % Improvement | .54 | .604 |
| Ryan, 1976 ¹⁸ | 12 | 12 | 24 | F statistic = 8.26 | .65 | .775 |
| Ryan, 1977 ¹⁷ | 10 | 10 | 20 | P = < .001 Improvement | .55 | .618 |
| Sausa, 1975 ⁵⁴ | 40 | 40 | 80 | % Improvement, TMAS | .35 | .365 |
| | | | | % Improvement, HAS | .45 | .485 |
| | | | | % Improvement, CRS | .35 | .365 |
| Schmitt, 1986 ²⁶ | 30 | 30 | 60 | P = < .05 STAI, State Anxiety | .35 | .365 |
| | | | | P = < .05 STAI, Trait Anxiety | .35 | .365 |
| | | | | P = < .05 IPAT | .35 | .365 |
| | | | | P = < .05 POMS Anxiety | .35 | .365 |
| Smith, 1999 ⁵⁵ | 23 | | 23 | t-score, State Anxiety | .74 | .950 |
| | | | | t-score, Trait Anxiety | .81 | 1.127 |
| Smith, 1975 ²⁷ | 36 | 36 | 72 | P = < .001, POMS Anxiety | .46 | .497 |
| Smith, 1992 ⁵⁶ | 31 | | 31 | % Improvement | .41 | .436 |
| Smith, 1994 ²² | 10 | 11 | 21 | P = < .05, POMS Anxiety | .60 | .693 |
| Smith, 2002 ⁵⁷ | 146 | 107 | 253 | P = < .03, POMS Anxiety | .19 | .192 |
| Taylor, 1991 ⁵⁸ | 15 | 15 | 30 | P = < .05, Diastolic BP | .50 | .549 |
| | | | | P = < .05, Systolic, BP | .50 | .549 |
| | | | | P = < .05, STAI, State | .50 | .549 |
| | | | | P = < .05, pulse rate | .50 | .549 |
| Von Richthoven, 1980 ⁵⁹ | 5 | 5 | 10 | P = < .001, Psychiatric Rating | .97 | 2.092 |
| | | | | P = < .005, STAI, State | .92 | 1.589 |
| | | | | P = < .005, Self Rating | .92 | 1.589 |
| Voris, 1995 ⁷ | 40 | 65 | 105 | P = < .0001, STAI, State | .49 | .536 |
| | | | | % Improvement, EMG | .63 | .741 |
| | | | | P = < .01, Temperature Change | .40 | .424 |
| Voris, 1996 ⁶⁰ | 8 | 7 | 15 | P = < .01, STAI, Trait | .80 | 1.099 |
| | | | | % Improvement, EMG | .53 | .590 |
| Weingarten, 1981 ⁶¹ | 12 | 12 | 24 | P = < .05, POMS Anxiety | .55 | .618 |
| Winick, 1999 ⁸ | 16 | 17 | 33 | P = < .02, Dentist RS | .56 | .633 |
| | | | | P = < .02, Patient SRS | .56 | .633 |

*Beck AI is the Beck Anxiety Index; Hamilton AS is the Hamilton Anxiety Scale, also known as HAS or HAMA; EMG is the electromyogram; STAI is the State/Trait Anxiety Inventory; TMAS is the Taylor Manifest Anxiety Scale; EEG is the electroencephalograph; State Anxiety and Trait Anxiety are both from the STAI; POMS is the Profile of Mood States; MACL is the Modified Adjective Check List; 16PF is the 16 Personality Factor Scale; CRS is a Clinical Rating Scale; IPAT is the International Personality and Ability Test; BP is blood pressure; RS is rating scale; SRS is self rated response scale

TABLE 6. LIST AND DESCRIPTION OF ANXIETY STUDIES

| AUTHOR | DIAGNOSIS | BLINDING | | | STUDY DESIGN | OUTCOME MEASURE |
|----------------------|--|----------|-----------|----------|------------------------------------|---|
| | | PATIENT | THERAPIST | ASSESSOR | | |
| Bianco, 1994 | Polysubstance Abusers | Yes | Yes | Yes | Double-blind | Beck/Hamilton Anxiety Scale |
| Feighner, 1973 | Psychiatric Inpatients | Yes | Yes | Yes | Double-blind Crossover | Global Rating Scale |
| Flemenbaum, 1974 | Psychiatric Outpatients | No | No | No | Open Clinical, Historical Controls | Global Rating Scale |
| Frankel, 1973 | Insomniacs | Yes | Yes | Yes | Double-blind, Crossover | TMAS |
| Gibson, 1983 | Outpatient Psychiatric | Yes | Yes | Yes | Double-blind | EMG, STAI |
| Gomez, 1974 | Heroin Addicts | Yes | Yes | Yes | Double-blind | TMAS |
| Hearst, 1974 | Outpatient Psychiatric | Yes | Yes | Yes | Double-blind | Self Ratings |
| Heffernan, 1996 | Outpatient Pain Patients | Yes | Yes | Yes | Double-blind | 4 Physiologic Measures |
| Heffernan, 1996a | Outpatient Pain Patients | Yes | Yes | Yes | Double-blind | EEG |
| Jamelka, 1975 | Prisoners, Psychiatric Ward | Yes | Yes | Yes | Double-blind | Hamilton AS |
| Kirsch, 2002 | Physicians' Report of Patient Response | No | No | No | Physician Survey | Physicians' Clinical Ratings |
| Krupitsky, 1991 | Alcoholic Inpatients | Yes | Yes | Yes | Double-blind | STAI, TMAS |
| Levitt, 1975 | Psychiatric Inpatients | Yes | Yes | Yes | Double-blind | TMAS |
| McKenzie, 1976 | Psychiatric Outpatients | No | No | No | Open Clinical | Skin Potential |
| Magora, 1967 | Psychiatric Inpatients | No | No | No | Open Clinical | Physician Clinical Rating |
| Matteson, 1986 | Graduate Students, Business School | No | No | No | Open Clinical | STAI |
| May, 1993 | Inpatient Drug Treatment | No | No | No | Open clinical | MAACL |
| Moore, 1975 | Outpatient Psychiatry | Yes | No | No | Crossover | Psychiatrist Ratings |
| Overcash, 1999 | Outpatient Psychiatry | No | No | No | Open Clinical | Physiological Measures, Self Rating Scale |
| Overcash, 1989 | Marijuana Patients | No | No | No | Open Clinical/different therapies | EMG, 16PF |
| Passini, 1976 | Inpatient Psychiatric | Yes | Yes | Yes | Double-blind | MAACL, STAI |
| Patterson, 1984 | Polydrug Abusers | No | No | No | Open Clinical | Abstinence Syndrome |
| Philip, 1991 | Polydrug Withdrawal | Yes | Yes | Yes | Double-blind | Visual Analog Scale |
| Rosenthal, 1972 | Psychiatric Outpatients | Yes | Yes | Yes | Double-blind | Psychiatrist Ratings |
| Rosenthal, 1970 | Psychiatric Outpatients | No | No | No | Open Clinical | Psychiatrist Ratings |
| Rosenthal, 1970a | Psychiatric Outpatients | No | No | No | Open Clinical | Psychiatrist Ratings |
| Ryan, 55 | Psychiatric Inpatients | Yes | Yes | Yes | Double-blind | STAI-State |
| Ryan, 1977 | Psychiatric Inpatients | Yes | Yes | Yes | Double-blind | STAI-State |
| Sousa, 1975 | Psychiatric Outpatients | Yes | Yes | Yes | Double-blind | TMAS, HAS, Clinical Rating Scale |
| Schmitt, 1986 | Inpatient Polydrug | Yes | Yes | Yes | Double-blind | POMS, IPAT, STAI |
| Smith, 1999 | Outpatient Psychiatry | No | No | No | Open Clinical | STAI |
| Smith, 1975 | Inpatient Addiction | Yes | No | Yes | Single Blind | POMS |
| Smith, 1992 | Outpatient Phobic | No | No | No | Open Clinical | Self Rating Scale |
| Smith, 1994 | Closed Head Injured | Yes | Yes | Yes | Double-blind | POMS |
| Smith, 2002 | Inpatient Polydrug | No | No | No | Retrospective | POMS |
| Taylor, 1991 | Normal volunteers | Yes | Yes | Yes | Double-blind | BP, Pulse Rate, STAI |
| Von Richthofen, 1980 | Anxiety Neurosis | Yes | Yes | Yes | Double-blind, Crossover | Psychiatrist RS, Self RS, STAI |
| Voris, 1995 | Prison Parolees, Sex Offenders | Yes | Yes | Yes | Double-blind | STAI, EMG, Temperature |
| Voris, 1996 | Prison Parolees, Sex Offenders | No | No | No | Open Clinical | STAI, EMG |
| Weingarten, 1981 | Inpatient Alcoholics | Yes | Yes | Yes | Double-blind | POMS |
| Winick, 1999 | Dental Patients | Yes | Yes | Yes | Double-blind | VAS |

$r=.44$ to $r=.70$ range.

There are numerous statistical considerations that must be taken into account in performing meta-analysis and Appendix B illustrates the most important ones. The non-statistician may find it useful to consider these factors and gain personal confidence in this valuable technique.

Conclusion

There have now been roughly 50 years of experience in the U.S. using CES as a non-pharmaceutical treatment for anxiety although it has yet to achieve ubiquitous status as a therapeutic modality. This is most likely due to the fact that few U.S. medical schools teach CES treatment as part of their curricula, and none of the seven or eight CES companies in U.S. history have had sufficient staff to visit physicians' offices in the ubiquitous manner of today's pharmaceutical representatives.

Yet when physicians who use or prescribe CES are asked about its effectiveness, they are generally enthusiastic, as are the majority of CES patients themselves. Patient response on surveys are even more significant because some CES device distributors have a 30 day period during which a patient can return the device at little or no cost if it proves ineffective. Less than 2% of patients return the devices for this reason, and almost none are returned by patients who use them in the suggested manner for the treatment of their anxiety (e.g., 20 minutes to one hour a day for the first three weeks, then as needed to prevent symptoms from returning). The fact that such devices can cost over \$1,000 makes the tendency to keep them even more impressive.

It is also noteworthy that among the more than 6,000 patients who have been involved in CES studies in the U.S., and from the thousands of patients who completed surveys, there have been no significant, negative side effects reported from the use of CES. The National Research Council evaluated the safety of CES for the FDA stating that, "...significant side effects or complications attributable to the procedure are virtually nonexistent."⁶⁷

From the data available, one would assume that CES will continue to receive greater attention from clinicians as more become aware of the safety and efficacy of this treatment for anxiety and the myriad of anxiety related disorders, especially chronic pain. ■

Appendix A. Example of Meta-Analysis Probability Conversions

For example, if percent improvement is reported, that percent figure is converted directly into the effect size, r . Similarly, Z scores are converted to r by the formula $r = \frac{Z}{\sqrt{N}}$.

If student t scores are given, they are converted into r by the formula $r = \sqrt{\frac{t^2}{t^2 + df}}$ where df is the degrees of freedom. If the author only gives the resulting probability figure, such as .05 or .01, one can convert that into the t score from published probability tables and compute r from the formula given just above.

When non-parametric statistics are reported, such as chi squared (χ^2), the author ordinarily reports the probability estimate obtained (.05, .01, etc.). In these cases, the r can also be obtained by converting the probability estimate into a t score.

There are other considerations that a diligent statistician must keep in mind when conducting meta-analysis. For example, an author might report finger temperature as a physiologic correlate of anxiety and report that the patients' average finger temperature rose from 91 to 94 degrees fahrenheit.⁷ For a clinical researcher in the field of biofeedback, that is a dramatic change, but how can it best be added to a meta-analysis? As important as it appears to a biofeedback therapist, it is in fact only a 3.3% improvement. That percent gain could be translated directly to an effect size r of .03 which would make it appear insignificant.

On the other hand, if a t score of 2.62 was derived from the patients before and after treatment data, from the formula given above (10 patients were treated, giving a $df = 9$) one can derive an effect size of $r=.66$. What makes that difference possible, and which effect size is the correct one? It is well known that temperature does not exist on a scale of one to 100 in humans. That is, it is not a 100 point scale. Therefore raw temperature scores must be adjusted accordingly.

One way to do that is to determine the criteria for the temperature range in humans. If the finger temperature range in subjects who would be well enough to be able to walk into a clinic to participate in a study is 95 to 101, or six points, then each temperature shift would be equal to 16.67 points on a 100 point scale and an r from this example, derived that way, would be .50. That is greater than the r of .03, but not as great as the r of .66.

In our example the r of .66 was derived from t scores, because the researcher had utilized actual temperature scores from the subjects who were in the study and compared both the subjects' actual mean finger temperature scores and the variance of all of those scores around the mean in arriving at the t score. That indicates that the range was not evenly divided from 95 to 101 among the research subjects, and was obviously much narrower if a change in three temperature points resulted in the 66% gain.

Therefore the total possible range of finger temperatures in normal people walking the streets is less important to the meta-analysis than those found among pain patients who are anxious. Pain related anxiety is known to restrict the finger temperature range considerably. That is what was found in the study used in this example, and that is why it was examined.

In one study, pre- and post-temperatures were given with no other information. Accordingly, the derived 8% improvement had to either be deleted from the meta-analysis or statistically dealt with after a very close reading of the original publication, since simply adding that r of .08 into the analysis would not only be in error but would unnecessarily skew or bias the results.

To summarize the problem, if a given data point measured does not ordinarily fall along a 100 point scale of variation, the percent change has to be adjusted for consistency before the number can be added to the analysis. T scores, F scores, probability scores, χ^2 and the like are the calculations used to determine the actual range.

Appendix B. Statistical Considerations in CES Meta-Analysis

Tabulation of CES Studies

The first step in a meta-analysis is locating all the available studies. The more inclusive the analysis, the greater the confidence that can be placed in its results. There are some typical, but not always obvious difficulties in finding all studies done on a subject. On one hand, the analyst may not have found many of the published studies due to simply not looking carefully enough. Some studies may have been published in a foreign language that the analyst can not read. Also, there may have been some very important studies that had very robust results, but were not published. A manufacturer of a product being used in the studies might be of help, if it is still in business and maintains a library of all the research in its field.

Another important problem is that many journals only publish positive studies and will routinely decline to publish studies with negative findings. That is not a problem here as there are journals that have published negative CES studies. Curiously, some of these same publications have then refused to publish CES studies with positive findings.

The meta-analysis has succeeded if the majority of studies that are available have been identified, since it can easily be shown that unless the analyst was specifically biased in the sampling procedure, the outcome would rarely be changed significantly if the remainder of the completed studies, and perhaps many more, appeared in the analysis. This point will be made clearer when the standard error of the mean and the confidence limits of the effect size are discussed later.

In analyzing the research in CES, extensive ground work was readily available in the form of a published comprehensive annotated bibliography listing all CES studies.²⁸ Table 5 shows the list of CES anxiety studies that can be subjected to meta-analysis, and Table 6 provides additional information about these studies.

Note the variety of statistics reported in the various studies in Table 5. In four of the studies, the authors only reported the percent of patients who improved, or in the case of Hearst, the percent of patients who no longer showed symptoms following the study.^{38,41,46,48} If no other information is available, these studies can not be converted into effect sizes, which refer not to the percent of patients who improved, but the average percent of improvement of the patients in the group.

If those four studies are removed and the meta-analysis is completed with the remaining studies there are a total of 67 data points and the meta-analysis yields an effect size of $r=.58$. Thus meta-analysis of these studies reveal a strong effect size for CES in the treatment of anxiety.

The Problem with Crossover Designs in CES Studies

Another problem is with studies designed to have patients cross over from an active treatment to a sham treatment group. It has been known since the 1950s that once treated with CES, patients tend to continue improving.⁶² This is why protocols calling for crossovers from treatment to sham groups are

counter-intuitive designs for a CES study. Because it has been shown that the group which previously had CES continues to improve, sometimes for months,⁶³ the control effect when these patients are crossed over into the sham treatment in the second half of the study is rendered ineffective and the study can not continue due to lack of effective controls. However, if the statistician can find data of the first part of the study prior to the crossover, the data can be utilized for meta-analysis. This distinction is not trivial and is important for any thoughtful clinician to ponder.

There are two crossover studies in Table 5 in which the sham controls ended up with the same outcomes as the treated patients, meaning that the study controls were "lost" (invalidated) in the crossover.^{38,59} We cannot put confidence in the treatment effects reported in those studies, thus they should rightfully be removed from the meta-analysis. Now, the more statistically valid revised Table 5 yields an effect size of $r=.54$.

Dealing with Multiple Measures in a Study

When an author publishes more than one measure of the treatment being studied, this must be taken into account as well. For example, in anxiety studies it is not unusual for a researcher to report results based upon patient self rating scales, clinician's rating scales, standardized anxiety tests, and physiological measures such as galvanic skin response, electroencephalograph changes, electromyography changes, and so forth. Some say the meta-analysis should assess all of that information, since it is all relevant to the effect CES has in reducing anxiety. Most statisticians disagree.

There is an argument that since there is error variance in any study (e.g., a patient's anxiety rating of 75 may contain input from an argument earlier in the day, the noise in the room, etc., rather than the typical level of anxiety in that patient over time), every measure from a given study that is included in the meta-analysis skews the final results along the unknown lines of that error bias. Most agree that only one effect size score should be included from any given study in order to reduce the contribution of each study's error bias to a minimum. But if only one variable is to be included, one has to have confidence that the statistician will select the right one for meta-analysis.

In 14 of the studies remaining in table 7.1 in which there is more than one variable reported, the statistician has a choice of which of the effect size estimates given to include (or exclude, depending on bias). Some statisticians will average two or more effect sizes presented in a given study and include that average in the analysis. However, *one of the strongest rules in statistics is that percentages can never be averaged*. To illustrate, in Table 5, Taylor has included measures of diastolic and systolic blood pressure and pulse rate, and scores from the State/Trait Anxiety Inventory.³⁸ Depending on the initial mean blood pressure measures in his study, a significant drop, say from 240 to 230, will mean something quite different than a

drop from 130 to 125, yet the effect size would be very similar. The change in pulse rate measure is obviously not affecting the same anxiety characteristics as the standardized psychological test of anxiety, which are measured on very different response scales.

Taylor must have assumed that each of the four measures utilized told us something different about anxiety in his patients and their response to CES. Since each of the gains were converted into a probability estimate, they are basically equated on the probability scale. If the four probability estimates that were provided are converted into Z_r scores, they can be averaged, which makes it possible to include them in the meta-analysis as a combined summary of the anxiety measures Taylor considered to be important. In so doing, a combined effect size of $r=.50$ was obtained from that study.

Why Correlations Can Not Be Averaged

The column for the Z_r score is an adjustment that is made due to the tendency of correlations, and therefore the r score distribution, to skew or distort trend and/or values as it reaches higher levels. For example if the mean anxiety improvement score of a group of patients is increased five points (from 15 to 20) as a result of treatment, that is a 33% improvement. If another group of anxiety patients also has their average score increased five points, from 60 to 65 that is only an 8% gain. Percent changes are not the same at the lower and upper ends of the percent distribution. By converting r scores into Z_r scores, they are normalized to a standard distribution for the purpose of combining them in the final effect size estimate without the bias inherent in skewing. The average Z_r score obtained can then be converted back into an r score by looking it up in published tables.

If the remaining studies in Table 5 are now analyzed with only one measure used per study 38 data points remain bringing the effect size of CES research on anxiety to $r=.57$. If the statistician had yielded to bias and included only the strongest effect size from each study, the effect size would be $r=.62$.

Even though it is now known that the type of study design has a limited effect on the outcome of meta-analyses,⁶⁴ there is still a widely held belief that double-blind studies are the gold standard of research, so any real truth about the effectiveness of a given type of treatment must come from them. It can be seen in Table 5 that there are 17 double-blind studies (not counting the double-blind crossover studies for reasons given above), and if those are subjected to meta-analysis, $r=.52$ is obtained. That is not impressively different from the $r=.57$ obtained when all the studies were analyzed.

Big Studies vs. Little Studies

There is a final consideration that often concerns readers that pertains to the number of subjects in a study. Isn't a study with 500 people better than one with only 10? Statisticians do not concern themselves about such things as much as clinicians who are schooled in pharmacology might. The statistical analysis will be penalized by the statistical tables to a lesser or greater extent depending on the number of people in the

study. For example, if there are only 10 people in a study there must be a t score of 2.26 to reach the .05 level of significance, whereas if there are 500 people in the study there need only be a t score of 1.97. Such t scores are made up of differences in the means of two groups of scores; but take into account the amount of variation of the scores around those means, which tends to rise with increasing numbers of subjects. A t of 1.97 can be easier to achieve with a larger number of subjects, everything else being equal, than a larger t . On the other hand, a significance level of .05 in either study means exactly the same thing: there is a 5% chance that the outcome occurred by chance alone, given the number of people selected for the study. Conversely that also means that the differences obtained between those study groups will probably be found in 95 out of every 100 similar studies with an identical number of subjects.

If all the effect sizes reported above are re-examined, and the effect size obtained is corrected for the number of subjects in each study, giving larger weight to those studies with more subjects, the $r=.58$ obtained for the total group remains $r=.58$.

The $r=.56$ obtained after deleting the crossover studies becomes $r=.63$ when corrected for the number of subjects in each study.

When all the studies were analyzed, but only one measure from each was included, an $r=.57$ was obtained, but that changed to an $r=.65$ when corrected for number of subjects in each study. When analyzing only the largest score among several presented, $r=.62$ was found, and when weight corrected, it rose to $r=.67$.

The double-blind studies yielded an $r=.53$ and that fell to $r=.52$ when adjusted for number of subjects in each study. The interesting thing about the analysis of the 17 double-blind studies is that the standard deviation (how the separate scores distributed around the mean) of the effect size scores was only .11, meaning that all the effect sizes from the 17 studies were closely clustered about the mean. Thus, while remaining the benchmark of good science, rigorous double-blind studies tend not to achieve as robust results as studies that allow for the addition of clinical skills and adjustments. This effect can be seen in some of the newer CES studies that are now conducting crossovers to open clinical trials where the clinician or patient is able to adjust the stimulation parameters more closely approximating normal usage (e.g., a higher current or more treatment time). This means that the crossover group tends to do better than the original blinded treatment group that may have been given less current for a shorter time.^{65,66}

Study Outcome Variation: the Standard Error of the Mean and Confidence Intervals

The standard deviation of the first analysis in which all score points were used was .36, meaning that the effect size scores varied widely about the mean. That was to be expected when including such diverse measures as electromyograms, pulse rates and scores on psychiatric rating scales, etc. The stan-

dard deviation combining all the studies, but using only one data point from each was .32, still representing a wide scattering of scores about the mean.

The standard error of the mean can be derived from standard deviations. Standard error of the mean is an indication of the limits in which the effect size would be expected to fall if another group of 20 studies were analyzed later.

If physicians knew in advance what measure of anxiety they intended to use in practice (e.g., the State/Trait Anxiety Inventory), then the statistician might be asked to meta-analyze just the studies which had used that measure. This would provide a closer look at what might be expected from using CES in practice. For example, in Table 5 ten studies reported either state ($N=9$) and/or trait ($N=5$) anxiety from the STAI. A meta-analysis of the state anxiety studies provides a strong effect size of $r=.60$ with a standard deviation of only .18. On the other hand, the studies of chronic, trait anxiety yield an effect size of $r=.68$ but the standard deviation rises to .33.

One last set of questions remain; might a very different effect size for CES for the treatment of anxiety result if one waited for another 20 studies to be published? Or better yet, might it be best to wait until another 60 studies are published following that additional 20. Would this provide greater certainty that CES is an effective treatment for anxiety?

To answer that, the statistician can utilize the standard error

of the mean effect size, that gives the confidence interval of the effect size obtained. Using the formula, where ES is the obtained effect size and N is the number of studies meta-analyzed, the result is the standard error of the mean effect size. *By definition, if the standard error of the mean is multiplied by 1.98, a range (+ or -) is obtained within which 95% of the effect sizes in subsequent meta-analyses will probably fall. If it is then multiplied by 2.63, the range is obtained within which 99% of scores are expected to fall.* As an example, from the data on the 17 double-blind studies, the effect size was $r=.57$, with a standard deviation of .11. That would yield a standard error of the mean effect size of .03. Multiply the .03 by 2.63 and an expected effect size of from $r=.59$ to $r=.65$ (the confidence interval) in 99 of the next 100 meta-analyses of groups of double-blinded CES studies of anxiety is obtained.

When the more typical situation is used, as above, in which results of 38 studies were carefully screened and allowed only one data point per study (using the average of the data points reported per study), a more likely scenario emerges for future meta-analysis on the treatment of anxiety with CES. That effect size was $r=.57$ and the standard error of the mean was .05. If that is multiplied by 2.63 it gives a range of from $r=.44$ to $r=.70$ for the expected effect size to be found in 99 out of 100 similar meta-analyses of groups of such studies in the future. ■

Daniel L. Kirsch, PhD, DAAPM, FAIS is an internationally renowned authority on electromedicine with 34 years of experience in the electromedical field. He is a board-certified Diplomate of the American Academy of Pain Management, Fellow of the American Institute of Stress, Member of the International Society of Neuronal Regulation, and a Member of Inter-Pain (an association of pain management specialists in Germany and Switzerland). He served as Clinical Director of The Center for Pain and Stress-Related Disorders at Columbia-Presbyterian Medical Center, New York City, and of The Sports Medicine Group, Santa Monica, California. Dr. Kirsch is the author of two books on CES titled, The Science Behind Cranial Electrotherapy Stimulation, 2nd Ed. published by Medical Scope Publishing Corporation, Edmonton, Alberta, Canada in 2002; and Schmerzen lindern ohne Chemie CES, die Revolution in der Schmerztherapie, Internationale Ärztegesellschaft für Energiemedizin, Austria 2000, in German. Best known for designing the Alpha-Stim CES and MET line of medical devices, Dr. Kirsch is Chairman of Electromedical Products International, Inc. of Mineral Wells, Texas, USA with additional offices in Europe and Asia. Dr. Kirsch can be reached at dan@epii.com.

Marshall F. Gilula, M.D. is a Diplomate of the American Board of Psychiatry and Neurology and a Diplomate of the American Board of Medical Electroencephalography. He is also a board-certified Instructor in Biofeedback and Neurotherapy (NBCB). In 1978 he was a US-USSR NIMH Exchange Scientist working with cranial electrotherapy stimulation and general psychophysiology techniques at the P.K. Anokhin Institute, Soviet Academy of Medical Sciences, Moscow. In 1983 Dr. Gilula was the first Motoyama-Ben Tōv Fellow at the Institute of Life Physics, Tokyo (Mitaka-shi), Japan and researched neuroelectric methodology and the EEG of altered states with Professor Hiroshi Motoyama. Dr. Gilula has had four years of residency and postdoctoral fellowship training in psychiatry and over seven years of postdoctoral training in neurology (neurophysiology and epilepsy). He has 40 years of experience in clinical psychiatry, and was in the Department of Neurology at the University of Miami School of Medicine from 1999 through 2003. Dr. Gilula was a Senior Fellow, Miami Center for Patient Safety, Department of Anesthesiology, University of Miami from 2003 through 2005. Dr. Gilula is President and CEO of the Life Energies Research Institute in Miami. He can be reached at mgilula@mindspring.com.

References

- (Continued from Part 1, March 2007 issue)
32. Rosenthal R. *Meta-analytic procedures for social research*. Newbury Park, California: Sage Publications. 1991. p 134.
 33. Wolf FM. *Meta-analysis; quantitative methods for research synthesis*. Newbury Park, California: Sage Publications. 1986. pp 31-33.
 34. Gilula MF and Kirsch DL. Cranial electrotherapy stimulation review: a safer alternative to psychopharmaceuticals in the treatment of depression. *Journal of Neurotherapy*. 2005. 9(2):7-26.
 35. Corroborated by two independent neurological presenters. *Annual Neurological Update*. Department of Neurology, University of Miami School of Medicine, Miami Beach, Florida. 2002.
 36. Feighner JP, Brown SL, and Olivier JE. Electro-sleep therapy: a controlled double-blind study. *Journal of Nervous and Mental Disease*. 1973. 157(2):121-128.
 37. Flierlbaum A. Cerebral electrotherapy (electro-sleep): an open clinical study with a six month follow-up. *Psychosomatics*. 1974. 15(1):20-24.
 38. Frankel BL, Buchbinder R, and Snyder F. Ineffectiveness of electro-sleep in chronic primary insomnia. *Archives of General Psychiatry*. 1973. 29:563-568.
 39. Hearst ED, Cloninger CR, Crews EL, and Cadoret RJ. Electro-sleep therapy: a double-blind trial. *Archives of General Psychiatry*. 1974. 30(4):463-466.
 40. Jamelka R. *Cerebral electrotherapy and anxiety reduction*. Master's Thesis, Stephen F. Austin State University, 1975.
 41. Krupitsky EM, Burakov AM, Karandoshova GF, Katsnelson J, Lebedev VP, Grinenko AJ, and Borodkin JS. The administration of transcranial electric treatment for affective disturbances therapy in alcoholic patients. *Drug and Alcohol Dependence*. 1991. 27:1-6.
 42. Levitt EA, James NM, and Flavell P. A clinical trial

- of electrosleep therapy with a psychiatric inpatient sample. *Australian and New Zealand Journal of Psychiatry*. 1975. 9(4):287-290.
43. McKenzie RE, Rosenthal SH, and Driessner JS. Some psycho-physiologic effects of electrical transcranial stimulation (electrosleep). In Wulfsohn, N.L. and Sances, A (Eds) *The Nervous System and Electric currents*. Plenum. New York. 1976. pp 163-167.
44. Magora F, Beller A, Assael MI, and Askenazi A. Some aspects of electrical sleep and its therapeutic value, in Wageneder, F.M. and St. Schuy (Eds.) *Electrotherapeutic Sleep and Electroanaesthesia*. Amsterdam:Excerpta Medica Foundation. 1967. International Congress Series No. 136. pp 129-135.
45. May B and May C. Pilot project using the Alpha-Stim 100 for drug and alcohol abuse. In Kirsch, D.L., *The Science Behind Cranial Electrotherapy Stimulation*. Medical Scope Publishing. Edmonton, Alberta, Canada. 2002. p 51.
46. Moore JA, Mellor CS, Standage KF, and Strong H. A double-blind study of electrosleep for anxiety and insomnia. *Biological Psychiatry*. 1975. 10(1):59-63.
47. Overcash SJ. A retrospective study to determine the efficacy of cranial electrotherapy stimulation (CES) on patients suffering from anxiety disorders. *Amer J of Electromedicine*. 1999. 16(1):49-51.
48. Overcash SJ and Siebenthal A. The effects of cranial electrotherapy stimulation and multisensory cognitive therapy on the personality and anxiety levels of substance abuse patients. *Amer J of Electromedicine*. 1989. 6(20):105-111.
49. Passini FG, Watson CG, and Herder J. The effects of cerebral electric therapy (electrosleep) on anxiety, depression, and hostility in psychiatric patients. *J of Nervous and Mental Disease*. 1976. 163(4):263-266.
50. Patterson MA, Firth J, and Gardiner R. Treatment of drug, alcohol and nicotine addiction by neuroelectric therapy: analysis of results over 7 years. *J of Bioelectricity*. 1984. 3(1&2):193-221.
51. Philip P, Demotes-Mainard J, Bourgeois M, and Vincent JD. Efficiency of transcranial electrostimulation on anxiety and insomnia symptoms during a washout period in depressed patients; a double-blind study. *Biological Psychiatry*. 1991. 29:451-456.
52. Rosenthal SH and Wulfsohn NL. Electrosleep: a preliminary communication. *J of Nervous and Mental Disease*. 1970. 151(2):146-151.
53. Rosenthal SH and Wulfsohn NL. Studies of electrosleep with active and simulated treatment. *Current Therapeutic Research*. 1970a. 12(3):126-130.
54. Sausa AD and Choudbury PC. A psychometric evaluation of electrosleep. *Indian J of Psychiatry*. 1975. 17:133-137.
55. Smith RB. Cranial electrotherapy stimulation in the treatment of stress related cognitive dysfunction, with an eighteen month follow up. *J of Cognitive Rehabilitation*. 1999. 17(6):14-18.
56. Smith RB and Shiromoto FN. The use of cranial electrotherapy stimulation to block fear perception in phobic patients. *Current Therapeutic Research*. 1992. 51(2):249-253.
57. Smith RB and Tyson R. The use of transcranial electrical stimulation in the treatment of cocaine and/or polysubstance abuse. In Kirsch, D.L., *The science behind cranial electrotherapy stimulation*. Medical Scope Publishing, Edmonton, Alberta, Canada: 2002. pp 68-69.
58. Taylor DN. *Effects of cranial transcutaneous electrical nerve stimulation in normal subjects at rest and during stress*. Ph.D. Dissertation, Brooklyn College of the City University of New York. 1991.
59. Von Richthofen CL and Mellor CS. Electrosleep therapy: a controlled study of its effects in anxiety neurosis. *Canadian J of Psychiatry*. 1980. 25(3):213-229.
60. Voris MD and Good S. Treating sexual offenders using cranial electrotherapy stimulation. *Medical Scope Monthly*. 1996. 3(11):14-18.
61. Weingarten E. The effect of cerebral electrostimulation on the frontalis electromyogram. *Biological Psychiatry*. 1981. 16(1):61-63.
62. Obrasow AN. Electrosleep therapy, in Licht, S (Ed.) *Therapeutic Electricity and Ultraviolet Radiation*. Vol 4. New Haven: Elizabeth Licht, 1959. Chapter 5.
63. Brotman P. Low-intensity transcranial electrostimulation improves the efficacy of thermal biofeedback and quieting reflex training in the treatment of classical migraine headache. *Amer J of Electromedicine*. 1989. 6(5):120-123, 1989. Ph.D. dissertation, City University Los Angeles. 1986. pp 1-117.
64. Rosenthal R. Meta-analytic procedures for social research. Sage Publications. Newbury Park, California. 1991. p 51.
65. Lichtbroun AS, Raicer MC, and Smith RB. The treatment of fibromyalgia with cranial electrotherapy stimulation. *J of Clinical Rheumatology*. 2001. 7(2):72-78.
66. Cork RC, Wood P, Ming N, Clifton S, James E, and Price L. The effect of cranial electrotherapy stimulation (CES) on pain associated with fibromyalgia. *The Internet J of Anesthesiology*. 2004. 8(2).
67. National Research Council, Division of Medical Sciences. *An evaluation of electroanesthesia and electrosleep*. FDA Contract 70-22, Task Order No. 20 (NTIS PB 241305). 1974. pp 1-54.