# **Cognadev Technical Report Series**



18<sup>th</sup> October, 2018

# The Binomial Effect Size Display (BESD)

#### Is this always an accurate index of effect?

I provide the definition, some warnings of the conditions under which it may not always produce accurate results, and some worked examples demonstrating those conditions. Overall, I think it's a pretty good 'quick approximation' ... but it is no substitute to having all the data at hand to calculate the actual effect/accuracy implied by a correlation/validity coefficient.



# The **BESD**

The Binomial Effect Size Display (BESD) was introduced in 1982 in an article:

Rosenthal, R., & Rubin, D.R. (1982). A simple, general purpose display of magnitude of experimental effect. *Journal of Educational Psychology*, 74, 2, 166-169.

From p. 166:

"The question addressed by BESD is What is the effect on the success rate (e.g., survival rate, cure rate, improvement rate, selection rate, etc.) of the institution of a certain treatment procedure? It displays the change in success rate (e.g., survival rate, cure rate, improvement rate, selection rate, etc.) attributable to a certain treatment procedure. "

From page 17 of Rosenthal, R., Rosnow, R.L., & Rubin, D.R. (2000) Contrasts and effect sizes in behavioral research: A correlational approach. Cambridge UK: Cambridge University Press. ISBN: 0-521-65980-9...

"Instead of concentrating on r<sup>2</sup>, we recommend using the point-biserial itself to create a display of the practical importance of the particular magnitude of effect (Rosenthal & Rubin, 1982). This is done simply by recasting r as a 2 x 2 contingency table, in which the rows correspond to the independent variable displayed as a table, in which the rows correspond to the independent variable displayed as a dichotomous predictor (e.g. experiment vs control) and the columns correspond to the dependent variable displayed as a dichotomous outcome (e.g. improved vs not-improved). The correlation between these two dichotomous variables is set to equal the obtained point-biserial r. The specific question addressed by this binomial effect size display (BESD) is: What is the effect on the success rate of the implementation of a certain procedure?

Table 2.4 {below} illustrates the BESD based on an r of .32, which was reported to be the average size of the effect of psychotherapy in an early report of a meta-analysis (Glass, 1976). To find the psychotherapy success rate of 66%, we computed .50 + r/2, and to find the control success rate of 34%, we computed .50 - r/2. In other words, r = .32 is equivalent to increasing the success rate from 34% to 66% (which in another case might mean, for example, reducing an illness rate or a death rate from 66% to 34%). Notice that the difference between the rate of improvement in the psychotherapy group and that in the control group (i.e., 66% - 34% = 32%) corresponds to the value of r times 100. These percentages should not, of course, be mistaken for the raw percentages in the actual data, but they can be interpreted as "standardized" percentages in order for all the margins to be equal. Another way of saying this is that an r of .32 (or an  $r^2$  of .10) will amount to a difference between rates of improvement of 34% and 66% if half the population received psychotherapy and half did not, and if half the population improved and half did not. "

<b>TABLE 2.4</b> BESD for $r = .32$							
	Treatme						
Condition	Improved	Not improved	Totals				
Psychotherapy	66	34	100				
Control	34	66	100				
Totals	100	100	200				
Control Totals	34 100	66 100					

In short, the BESD is the probability of a random outcome (0.5) plus one-half of a *point-biserial* correlation coefficient (which is mathematically equivalent to a Pearson correlation and the conventional phi coefficient calculated from 2x 2 contingency tables).

But, as Hsu, L.M. (2004) Biases of success rate differences shown in Binomial Effect Size Displays. *Psychological Methods*, 9, 2, 183-197 points out ..

"**Abstract:** The intent of a binomial effect size display (BESD) is to show "the [real-world] importance of [an] effect indexed by a correlation [r]" (R. Rosenthal, 1994, p. 242) by reexpressing this correlation as a success rate difference (SRD) (e.g., treatment group success rate - control group success rate). However, SRDs displayed in BESDs generally overestimate realworld SRDs implied by correlations of (a) dichotomous X and Y variables ( $\phi$  coefficients), (b) dichotomous X and continuous Y variables (point-biserial coefficients [r<sub>pb</sub>s]), and (c) continuous X and Y variables (r<sub>xy</sub>s). Furthermore, overestimation biases are larger for r<sub>xy</sub>s than for r<sub>pb</sub>s. Differences in the sizes of biases linked to different correlations suggest that BESD SRDs reported for different correlations are not comparable. The stochastic difference index (N. Cliff, 1993; A. Vargha & H. D. Delaney, 2000) is recommended as an alternative to the BESD."

As Hsu says on p. 183 ..

"What may not be apparent from Rosenthal and Rubin's illustration is that the equality of the BESD SRD (calculated from Equation 1) and actual SRD of a 2 x 2 table does not generalize to 2 x 2 tables that do not have uniform marginal distributions."

Eq. 1 is: 
$$BESD \_ SRD = \left[ \left( .5 + \frac{r}{2} \right) - \left( .5 - \frac{r}{2} \right) \right] = r$$

And this gentle statement on page 18 from Rosenthal, R., Rosnow, R.L., & Rubin, D.R. (2000) Contrasts and effect sizes in behavioral research: A correlational approach. Cambridge UK: Cambridge University Press.

"It can be shown that the BESD is most appropriate when the variances within the two conditions are similar, as they are assumed to be whenever we compute the usual t statistics and its associated p-value."

Ok - all this is fine and dandy as it goes - but it's still a bit disconnected from what we really want to use the BESD for .. converting a correlation computed from a 2 x 2 decision-table to an index of classification accuracy. So, if we observe a correlation of 0.32 between a predictor (dichotomized into high-low say using a cut-score) and binary outcome (success, failure), then we can express that

correlation as how well we can classify outcomes above a chance level of 50%. That chance level is actually no more than a base-rate of positive outcome of .5 (50%).

For example, if we observe a correlation between those scoring above 70 and below 70 on a competency, and rated job success after 1 year, of .4, our overall classification accuracy would be

$$BESD = \left(.5 + \frac{r}{2}\right) = \left(.5 + \frac{.4}{2}\right) = .7 \text{ or } 70\% \text{ classification accuracy.}$$

So, let's see a few 'full feature' worked examples to really get to grips with the BESD, warts and all. I'm using my Dichot 3.1 software – which is freely available for download from the web: http://www.pbarrett.net/Dichot3/Dichot3.html

and includes a pdf manual and web-help embodies in the program.

#### The Rosenthal and Rubin example presented in their Table 2.4 (on my page 1 above) ...



For an explanation & formulae for all the other coefficients reported here – download the program manual (if you haven't already downloaded the software). <u>http://www.pbarrett.net/Dichot3/Dichot3\_1\_program\_manual.pdf</u>

# Where the Base Rate remains at 0.5 (50%), but the error-rates are no longer balanced (unequal marginals)

DIGHOT V.3.1 - Dicnotomous Relationship	s and Decision	1 Table Statistic	.5		
2	ARIABLE 1 (A	ctual/Disease,	Outcome)	Generate a Report	Compute
Yn Pres VARIABLE 2 (Readiated) Present/Success	1 es/Agree ent/Success 38 A	Expected Frequencies 21.5000	0 No/Disagree Expected Absent/Fail Frequencies 5 B 21.5000	MARGINALS 43	
Factor/	e Positive (TP)		False Positive (FP)		
Treatment) 0 No/Disagree	62 <b>C</b>	78.5000	95 <b>D</b> 78.5000	157	
Fals	se Negative (FN	0	True Negative (TN)		
	100		100	200 TOTAL N	
					негр
Coefficient	Value	Probability		Coefficient	Value
Pearson Chi-Square	32.2619	0.0000000		Sensitivity	0.3800
Likelihood Ratio	35.6867	0.0000000		Sensitivity Quality	0.2102
Pearson r/Phi	0.4016	0.0000000		Specificity	0.9500
Phi/Phi-Max	0.7674			Specificity Quality	0.7674
Yule's Q (Gamma)	0.8418	0.0000000			
Jaccard	0.3619		Positive Po	ower to Predict (PPP, PPV, PVP)	0.8837
G-Index (Hamman)	0.3300		Negative Po	wer to Predict (NPP, NPV, PVN)	0.6051
Bennett's B-Index	0.2662			Level (Q)	0.2150
Cohen's Kappa	0.3300			Overall Classification Assurance	0.6650 4
			E.	loo Dositivo Pato (Falso Alarm)	0.0000
Cohen d' Effect size	1.3413		Fa	False Negative Rate	0.0300
	0.5570			Faise Negative Rate	0.0200
Estimated Pearson r from d'			Odds of Outcom	ne given treatment (predicted)	7.6000
Estimated Pearson r from d'					
Estimated Pearson r from d' Base Rate	0.5000		Odds of Outcome NOT give	en treatment (or not predicted)	0.6526
Estimated Pearson r from d' Base Rate	0.5000		Odds of Outcome NOT give	en treatment (or not predicted) Odds Ratio	0.6526 11.6452

The  $BESD = \left(.5 + \frac{.4016}{2}\right) = .7008$  or 70.08% effect display is <u>not</u> equal to the overall classification accuracy of 66.50%.

#### DICHOT v.3.1 - Dichotomous Relationships and Decision Table Statistics **Generate a Report** Compute VARIABLE 1 (Actual/Disease/Outcome) 0 Yes/Agree No/Disagree Expected Expected MARGINALS Absent/Fail Present/Success Frequencies Frequencies Yes/Agree VARIABLE 2 Α в 40 36 20.0000 20.0000 4 Present/Success (Predicted/ False Positive (FP) True Positive (TP) Factor/ No/Disagree Treatment) 0 С D 80.0000 160 64 80.0000 96 Absent/Fail False Negative (FN) True Negative (TN) 100 100 TOTAL N 200 Help Coefficient Value Probability Coefficient Value 32.0000 Sensitivity 0.3600 Pearson Chi-Square 0.0000000 Sensitivity Quality 0.2000 **Likelihood Ratio** 35.8885 0.0000000 Specificity 0.9600 Pearson r/Phi 0.4000 0.0000000 **Specificity Quality** 0.8000 Phi/Phi-Max 0.8000 Yule's Q (Gamma) 0.8621 0.0000000 Positive Power to Predict (PPP, PPV, PVP) 0.9000 Jaccard 0.3462 Negative Power to Predict (NPP, NPV, PVN) 0.6000 G-Index (Hamman) 0.3200 Level (Q) 0.2000 **Bennett's B-Index** 0.2487 Cohen's Kappa 0.3200 **Overall Classification Accuracy** 0.6600 False Positive Rate (False Alarm) 0.0400 Cohen d' Effect size 1.3946 **False Negative Rate** 0.6400 Estimated Pearson r from d' 0.5720 Odds of Outcome given treatment (predicted) 9.0000 0.5000 Odds of Outcome NOT given treatment (or not predicted) **Base Rate** 0.6667 **Odds Ratio** 13.5000 Relative Improvement Over Chance (RIOC) 1.6000 **Relative Risk** 2.2500

#### Where the Base Rate remains at 0.5 (50%), but the error-rates are no longer balanced (unequal marginals)

The  $BESD = \left(.5 + \frac{.4}{2}\right) = .7$  or 70.0% effect display is <u>not</u> equal to the overall classification accuracy (effect) of 66%.

# Where the Base Rate is just above chance 0.537 (54%), and the error-rates are not balanced (unequal marginals)

DICHOT v.3.1 - Dichotomous Relationships and Decision Table Statistics					• 🗙
	VARIABLE 1 ()	Actual/Disease/	Outcome) Generat	e a Report Com	pute
YARIABLE 2 (Predicted/ Factor/ Truesent/Success No/Disagree	1 /es/Agree sent/Success 36 A we Positive (TP)	Expected Frequencies <b>19.8704</b>	0 No/Disagree Expected Absent/Fail Frequencies MARGINALS 1 B 17.1296 37 False Positive (FP)		
Absent/Fail	80 C	96.1296	99 <b>D</b> 82.8704 179		
	116	.,	100 216	TOTAL N	Help
Coefficient	Value	Probability		Coefficient Valu	Je
Pearson Chi-Square	34.1269	0.0000000		Sensitivity 0.31	03
Likelihood Ratio	42.9326	0.0000000	Sensitiv	vity Quality 0.16	78
Pearson r/Phi	0.3975	0.0000000		Specificity 0.99	00
Phi/Phi-Max	0.9416		Specifi	city Quality 0.94	16
Yule's Q (Gamma)	0.9561	0.0000000			
Jaccard	0.3077		Positive Power to Predict (PPP	, PPV, PVP) 0.97	30
G-Index (Hamman)	0.2500		Negative Power to Predict (NPP,	NPV, PVN) 0.55	31
Bennett's B-Index	0.1781			Level (Q) 0.17	13
Cohen's Kappa	0.2848		Overall Classificatio		50
			False Positive Rate (F	alse Alarm) 0.01	100
Cohen d' Effect size	1.8355		False Positive Nate (F	gative Rate 0.65	897
Estimated Pearson r from d'	0.6752			Successfully 0.00	
			Odds of Outcome given treatment	(predicted) 36.00	000
Base Rate	0.5370		Odds of Outcome NOT given treatment (or not	predicted) 0.80	)81
				Odds Ratio 44.55	500
Relative Improvement Over Chance (RIOC)	3.0991		R	elative Risk 2.17	770

The  $BESD = \left(.5 + \frac{.3975}{2}\right) = .6988$  or 69.88% effect display is <u>not</u> equal to the overall classification accuracy (effect) of 62.5%.

### Where the Base Rate is substantively above chance 0.5941 (59%), and the error-rates are not balanced (unequal marginals)

DICHOT v.3.1 - Dichotomous Relationship	os and Decision Table	Statistics		- • •	
	VARIABLE 1 (Actual,	Disease/Outcome)	Generate a Report	Compute	
Y Pres VARIABLE 2 (Predicted/ Factor/ Treatment) 0 No/Disagree Absent/Fail Fail Fail	1 Exp   Yes/Agree Exp   sent/Success Freque   70 A 40   ue Positive (TP) 100 C 12   lse Negative (FN) 100 12	0 Expected   No/Disagree Expected   Absent/Fail Frequencies   5926 4 B 27.4074   False Positive (FP)   .4074 96 D 72.5926   True Negative (TN)	MARGINALS 74 196		
	170	100	270 TOTAL N	Help	
Coefficient	Value Pro	ability	Coefficient	Value	I ne bottom line, for me, is that the
Pearson Chi-Square	43.7382 0.0	00000	Sensitivity	0.4118	BESD is a reasonable stab at what
Likelihood Ratio	53.1884 0.0	00000	Sensitivity Quality	0.1897	the overall classification accuracy
Pearson r/Phi	0.4025 0.0	00000	Specificity	0.9600	would be for a 2 x 2 decision-table
Phi/Phi-Max	0.8541		Specificity Quality	0.8541	
Yule's Q (Gamma)	0.8876 0.0	00000			phi (correlation/validity coefficient)
Jaccard	0.4023	Positive Po	ower to Predict (PPP, PPV, PVP)	0.9459	
G-Index (Hamman)	0.2296	Negative Pov	wer to Predict (NPP, NPV, PVN)	0.4898	But it can seriously mislead if the
Bennett's B-Index	0.2192		Level (Q)	0.2741	
Cohen's Kappa	0.3104		Overall Classification Accuracy	0.6148	error-rates are not balanced in tern
		Fa	Ise Positive Rate (False Alarm)	0.0400	of yielding near-equal marginal, as
Cohen d' Effect size	1.5295		False Negative Rate	0.5882	Hsu (2004) indicates
Estimated Pearson r from d'	0.5941				
		Odds of Outcom	e given treatment (predicted)	17.5000	BESD SKDS tend to overestimate
Base Rate	0.6296	Odds of Outcome NOT give	n treatment (or not predicted)	1.0417	targeted real-world SRDs in virtual
			Odds Ratio	16.8000	all real-world applications" (p. 195)
Relative Improvement Over Chance (RIOC)	-0.0118		Relative Risk	1.8541	(pr 199)

The  $BESD = \left(.5 + \frac{.4025}{2}\right) = .7013$  or 70.13% effect display is <u>not</u> equal to the overall classification accuracy (effect) of 61.48%.

9 | Page