

Implementing ODA from Within Stata: Exploratory Hypothesis, Binary Class Variable, Categorical Ordinal Attribute

Paul R. Yarnold, Ph.D. and Ariel Linden, Dr.P.H.
Optimal Data Analysis, LLC Linden Consulting Group, LLC

This paper describes how an exploratory (post hoc, nondirectional, two-tailed) hypothesis involving a binary (dichotomous) class variable and a categorical ordinal (three-level) attribute is evaluated using MegaODA software using the new Stata package implementing ODA analysis.

Recent papers¹⁻¹³ introduce the new Stata package called **oda**¹⁴ for implementing ODA from within the Stata environment. Because this package is a wrapper for the MegaODA software system¹⁵⁻¹⁷, the MegaODA.exe file must be loaded on the computer for the **oda** package to work (MegaODA software is available at <https://odajournal.com/resources/>). To download the **oda** package, at the Stata command line type: “ssc install oda” (without the quotation marks). This paper demonstrates use of the **oda** package to evaluate a nondirectional hypothesis involving a binary class variable, and an ordinal (three-level) attribute.

Methods

Data

We consider data from Snyder, Wills, and Grady-Fletcher comparing outcomes of two types of marital therapy for unhappily married couples.¹⁸ Arbitrary dummy-codes identified two types of *therapy*: insight=1, behavioral=2.

Therapy *outcome* was rated on a three-category ordinal scale ranging from worst to best outcome: divorced=1, no change=2, improved=3. Data for every subject was entered in free format on a separate line as space-delimited text (ASCII) characters.¹⁹

Analytic Process

We repeat the ODA analysis previously conducted on these data (see example 5.6, *Optimal Data Analysis: A Guidebook with Software for Windows*²⁰). The nondirectional or “two-tailed” alternative hypothesis is that the binary class (“dependent”) variable *therapy* can be discriminated on the basis of *outcome* (ordinal attribute or “independent variable”). The null hypothesis is that this is not true. Weighting by prior odds (the default setting) is used to obtain a model which maximizes ESS (i.e., classification accuracy normed vs. chance), and a total of 25,000 Monte Carlo iterations are used to estimate Type I error (i.e., *p* value).²⁰

For these data, **oda** is implemented using the following syntax to test the *a priori* hypothesis for the attribute *outcome* (see the **oda** help file for a complete description of syntax options):

```
oda therapy outcome, pathoda("C:\ODA\")
store("C:\ ODA\output") iter(25000)
```

The above syntax is explained as follows: The variable “therapy” is the *class* variable; the variable “outcome” is the *attribute*; the directory path where the MegaODA.exe file is located on the computer is “C:\ODA\”; the directory path where the output and other files generated during the analysis are stored is “C:\ODA\output”; and 25,000 iterations (repetitions) are used to compute the permutation *p*-value.

The **oda** package produces an extract of the total output produced by the ODA software (the complete output is stored in the specified directory with the extension “.out”).

As seen in the **oda** output, the ODA model is interpreted as follows: “if *outcome* ≤ 1.5 then predict *therapy* = 2; otherwise, predict *therapy* = 1.” Couples in behavioral therapy were predicted to be divorced, whereas couples in insight therapy were predicted to show no change or improvement. As seen, this model correctly classified 89.66% of the couples in behavioral therapy, but only 46.15% of the couples in insight therapy.

Effect strength for sensitivity (ESS) is labelled in the output as “Effect Strength PAC” (Percentage Accurate Classification). The ESS is 35.81% which corresponds to an effect of moderate strength²⁰ with permutation *p*<0.0055. In summary, ODA identified a model which discriminated therapeutic methods with moderate strength, and this effect was statistically significant.

```
ODA model:
-----
IF OUTCOME <= 1.5 THEN THERAPY = 2
IF 1.5 < OUTCOME THEN THERAPY = 1
```

Summary for Class THERAPY Attribute OUTCOME

Performance Index	Train
Overall Accuracy	69.09%
PAC THERAPY=1	89.66%
PAC THERAPY=2	46.15%
Effect Strength PAC	35.81%
PV THERAPY=1	65.00%
PV THERAPY=2	80.00%
Effect Strength PV	45.00%
Effect Strength Total	40.40%

Monte Carlo summary (Fisher randomization):

```
-----
Iterations: 25000
Estimated p: 0.005480
```

Discussion

This paper shows how to use ODA to identify the model that maximally discriminates between any two categories of a class variable using a categorical ordinal attribute.

ODA should be considered the preferred approach over other methods because it avoids statistical assumptions required of conventional models, is insensitive to skewed data or outliers, and has the ability to handle any variable metric including categorical, Likert-type integer, and real number measurement scales.²⁰ Moreover, in contrast to other methods, ODA also has the unique ability to ascertain optimal (maximum-accuracy) assignments (categorical attributes) or cutpoints (ordered attributes) on the attribute, which facilitates the use of measures of predictive accuracy. Furthermore, ODA can perform cross-validation using LOO (and many other methods²⁰) which allows for assessment of potential cross-generalizability of the model to independent random samples.

For these reasons we recommend that researchers employ ODA and CTA frameworks to evaluate the statistical hypotheses which are explored in their laboratory and field research endeavors.²¹⁻³⁵

References

- ¹Linden A (2020). Implementing ODA from within Stata: An application to data from a randomized controlled trial (*Invited*). *Optimal Data Analysis*, 9, 9-13.
- ²Linden A (2020). Implementing ODA from within Stata: Implementing ODA from within Stata: An application to estimating treatment effects using observational data (*Invited*). *Optimal Data Analysis*, 9, 14-20.
- ³Linden A (2020). Implementing ODA from within Stata: An application to dose-response relationships (*Invited*). *Optimal Data Analysis*, 9, 26-32.
- ⁴Linden A (2020). Implementing ODA from within Stata: assessing covariate balance in observational studies (*Invited*). *Optimal Data Analysis*, 9, 33-38.
- ⁵Linden A (2020). Implementing ODA from within Stata: Evaluating treatment effects for survival (time-to-event) outcomes (*Invited*). *Optimal Data Analysis*, 9, 39-44.
- ⁶Linden A (2020). Implementing ODA from within Stata: Evaluating treatment effects in multiple-group interrupted time series analysis (*Invited*). *Optimal Data Analysis*, 9, 45-50.
- ⁷Linden A (2020). Implementing ODA from within Stata: identifying structural breaks in single-group interrupted time series designs (*Invited*). *Optimal Data Analysis*, 9, 51-56.
- ⁸Linden A (2020). Implementing ODA from within Stata: Finding the optimal cut-point of a diagnostic test or index (*Invited*). *Optimal Data Analysis*, 9, 74-78.
- ⁹Yarnold PR, Linden A (2020). Implementing ODA from within Stata: Exploratory hypothesis, binary class variable, and binary attribute. *Optimal Data Analysis*, 9, 94-98.
- ¹⁰Yarnold PR, Linden A (2020). Implementing ODA from within Stata: Confirmatory hypothesis, binary class variable, and binary attribute. *Optimal Data Analysis*, 9, 99-103.
- ¹¹Yarnold PR, Linden A (2020). Implementing ODA from within Stata: Exploratory hypothesis, binary class variable, and binary attribute. *Optimal Data Analysis*, 9, 104-108.
- ¹²Yarnold PR, Linden A (2020). Implementing ODA from within Stata: Exploratory hypothesis, binary class variable, and ordinal (rank) attribute. *Optimal Data Analysis*, 9, 109-113.
- ¹³Yarnold PR, Linden A (2020). Implementing ODA from within Stata: confirmatory hypothesis, binary class variable, and ordinal attribute. *Optimal Data Analysis*, 9, 128-132.
- ¹⁴Linden A (2020). ODA: Stata module for conducting Optimal Discriminant Analysis. *Statistical Software Components S458728*, Boston College Department of Economics.
- ¹⁵Soltysik RC, Yarnold PR (2013). MegaODA large sample and BIG DATA time trials: Separating the chaff. *Optimal Data Analysis*, 2, 194-197.
- ¹⁶Soltysik RC, Yarnold PR (2013). MegaODA large sample and BIG DATA time trials: Harvesting the Wheat. *Optimal Data Analysis*, 2, 202-205.
- ¹⁷Yarnold PR, Soltysik RC (2013). MegaODA large sample and BIG DATA time trials: Maximum velocity analysis. *Optimal Data Analysis*, 2, 220-221.
- ¹⁸Snyder DK, Wills R, Grady-Fletcher A (1991). Long-term effectiveness of behavioral versus insight-oriented marital therapy: A four-year follow up study. *Journal of Consulting and Clinical Psychology*, 59, 138-146.

- ¹⁹Bryant FB, Harrison PR (2013). How to create an ASCII input data file for UniODA and CTA software (*Invited*). *Optimal Data Analysis*, 2, 2-6.
- ²⁰Yarnold PR, Soltysik RC (2005). *Optimal data analysis: Guidebook with software for Windows*. Washington, D.C.: APA Books.
- ²¹Rhodes NJ (2020). Statistical power analysis in ODA, CTA and Novometrics (*Invited*). *Optimal Data Analysis*, 9, 21-25.
- ²²Yarnold PR, Linden A. (2016). Novometric analysis with ordered class variables: The optimal alternative to linear regression analysis, *Optimal Data Analysis*, 5, 65-73.
- ²³Yarnold PR, Linden A (2016). Theoretical aspects of the D statistic. *Optimal Data Analysis*, 22, 171-174.
- ²⁴Linden A, Yarnold PR, Nallomothu BK (2016). Using machine learning to model dose-response relationships. *Journal of Evaluation in Clinical Practice*, 22, 860-867.
- ²⁵Linden A, Yarnold PR (2017). Using classification tree analysis to generate propensity score weights. *Journal of Evaluation in Clinical Practice*, 23, 703-712.
- ²⁶Linden A, Yarnold PR (2017). Modeling time-to-event (survival) data using classification tree analysis. *Journal of Evaluation in Clinical Practice*, 23, 1299-1308.
- ²⁷Linden A, Yarnold PR (2018). Identifying causal mechanisms in health care interventions using classification tree analysis. *Journal of Evaluation in Clinical Practice*, 24, 353-361.
- ²⁸Linden A, Yarnold PR (2017). Minimizing imbalances on patient characteristics between treatment groups in randomized trials using classification tree analysis. *Journal of Evaluation in Clinical Practice*, 23, 1309-1315.
- ²⁹Linden A, Yarnold PR (2018). Estimating causal effects for survival (time-to-event) outcomes by combining classification tree analysis and propensity score weighting. *Journal of Evaluation in Clinical Practice*, 24, 380-387.
- ³⁰Linden A, Yarnold PR (2016). Using machine learning to assess covariate balance in matching studies. *Journal of Evaluation in Clinical Practice*, 22, 848-854.
- ³¹Linden A, Yarnold PR (2016). Using machine learning to identify structural breaks in single-group interrupted time series designs. *Journal of Evaluation in Clinical Practice*, 22, 855-859.
- ³²Linden A, Yarnold PR (2016). Combining machine learning and matching techniques to improve causal inference in program evaluation. *Journal of Evaluation in Clinical Practice*, 22, 868-874.
- ³³Linden A, Yarnold PR (2016). Combining machine learning and propensity score weighting to estimate causal effects in multivalued treatments. *Journal of Evaluation in Clinical Practice*, 22, 875-885.
- ³⁴Linden A, Yarnold PR (2018). Using machine learning to evaluate treatment effects in multiple-group interrupted time series analysis. *Journal of Evaluation in Clinical Practice*, 24, 740-744.
- ³⁵Yarnold PR, Rhodes NJ, Linden A (2020). Selecting an appropriate weighting strategy in maximum-accuracy time-to-event (survival) analysis. *Optimal Data Analysis*, 9, 3-6.

Author Notes

No conflicts of interest were reported.