

# Inference on the Covariate-Specific Overlap Coefficient (OVL)

C.T. Nakas<sup>1,2</sup>, A.M. Franco-Pereira<sup>3</sup>, M.C. Pardo<sup>3</sup>, B. Reiser<sup>4</sup>

<sup>1</sup>University of Thessaly

<sup>2</sup>Inselspital/University of Bern

<sup>3</sup>Complutense University of Madrid

<sup>4</sup>University of Haifa

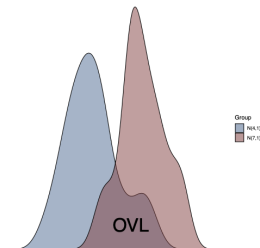
19<sup>th</sup> International Summer Conference on Probability and Statistics

# The Overlap Coefficient (OVL)

- Measures similarity between two distributions:

$$\text{OVL} = \int \min(f_1(x), f_2(x)) dx$$

- Ranges from 0 (no overlap) to 1 (identical)
- Common in diagnostic testing, bioequivalence, similarity metrics
- Not directional: unlike AUC, it captures general similarity



## Matusita / Weitzman Overlap Coefficient:

$$\text{OVL}(f_1, f_2) = \int \min(f_1(x), f_2(x)) dx$$

## Czekanowski / Sørensen-Dice Coefficient:

$$C(f_1, f_2) = \frac{2 \int \min(f_1(x), f_2(x)) dx}{\int f_1(x) dx + \int f_2(x) dx}$$

*Note: If  $f_1$  and  $f_2$  are PDFs (i.e., integrate to 1), then:*

$$C(f_1, f_2) = \text{OVL}(f_1, f_2)$$

# Morisita Index and Comparison

## Morisita Overlap Index:

$$C_M(f_1, f_2) = \frac{2 \int f_1(x) f_2(x) dx}{\int f_1^2(x) dx + \int f_2^2(x) dx}$$

## Comparison of Theoretical Overlap Indices:

Coefficient	Symmetric	Range	Equals OVL for PDFs?
OVL (Matusita/Weitzman)	Yes	[0, 1]	—
Czekanowski/Sørensen-Dice	Yes	[0, 1]	Yes
Morisita	Yes	[0, 1]	No

# Main reference #1

Franco-Pereira AM, Nakas CT, Reiser B, Carmen Pardo M. Inference on the overlap coefficient: The binormal approach and alternatives. Stat Methods Med Res. 2021 Dec;30(12):2672-2684. doi: 10.1177/09622802211046386. Epub 2021 Oct 23. PMID: 34693817.

**Sage Journals**   [Advanced search](#)

Browse by discipline  Information for

**Statistical Methods in Medical Research**

Impact Factor: **1.9** / 5-Year Impact Factor: **2.5**

Available access | Research article | First published online October 23, 2021

**Inference on the overlap coefficient: The binormal approach and alternatives**

[Alba Maria Franco-Pereira](#) [Christos T. Nakas](#) i.-l. and [Maria Carmen Pardo](#) [View all authors and affiliations](#)

[Volume 30, Issue 12](#) | <https://doi.org/10.1177/09622802211046386>

Contents | PDF/EPUB | Cite article | Share options | Information, rights and permissions

## Abstract

The overlap coefficient (*OVL*) measures the similarity between two distributions through the overlapping area of their distribution functions. Given its intuitive description and ease of visual representation by the straightforward depiction of the amount of overlap between the two corresponding histograms based on samples of measurements from each one of the two distributions, the development of accurate methods for confidence interval construction can be useful for applied researchers. The overlap coefficient has received scant attention in the literature since it lacks readily available software for its implementation, while inferential procedures that can cover the whole range of distributional scenarios for the two underlying distributions are missing. Such methods, both parametric and non-parametric are developed in this article, while R-code is provided for their implementation. Parametric approaches based on the binormal model show better performance and are appropriate for use in a wide range of distributional scenarios. Methods are assessed through a large simulation study and are illustrated using a dataset from a study on human immunodeficiency virus-related cognitive function assessment.

# Estimation of OVL

Let  $X_{11}, X_{12}, \dots, X_{1n_1}$  and  $X_{21}, X_{22}, \dots, X_{2n_2}$  drawn from independent  $N(\mu_1, \sigma_1^2)$  and  $N(\mu_2, \sigma_2^2)$  respectively.

$$OVL = 1 + \Phi\left(\frac{x_1 - \mu_1}{\sigma_1}\right) - \Phi\left(\frac{x_1 - \mu_2}{\sigma_2}\right) - \Phi\left(\frac{x_2 - \mu_1}{\sigma_1}\right) + \Phi\left(\frac{x_2 - \mu_2}{\sigma_2}\right)$$

assuming  $\sigma_1^2 < \sigma_2^2$ , where

$$x_1 = \frac{(\mu_1\sigma_2^2 - \mu_2\sigma_1^2) - \sigma_1\sigma_2\sqrt{(\mu_1 - \mu_2)^2 + (\sigma_1^2 - \sigma_2^2)\log\left(\frac{\sigma_1^2}{\sigma_2^2}\right)}}{(\sigma_2^2 - \sigma_1^2)}$$

$$x_2 = \frac{(\mu_1\sigma_2^2 - \mu_2\sigma_1^2) + \sigma_1\sigma_2\sqrt{(\mu_1 - \mu_2)^2 + (\sigma_1^2 - \sigma_2^2)\log\left(\frac{\sigma_1^2}{\sigma_2^2}\right)}}{(\sigma_2^2 - \sigma_1^2)}$$

where  $\Phi$  is the cumulative distribution function of a  $N(0, 1)$ .

Denote with  $\widehat{OV_L}$  when plugging-in the MLEs.

Since  $\hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}_1, \hat{\sigma}_2$  are all independent, using the delta method, we obtain

$$\text{Var}(\widehat{OV_L}) \approx \left(\frac{\partial OV_L}{\partial \mu_1}\right)^2 \text{Var}(\hat{\mu}_1) + \left(\frac{\partial OV_L}{\partial \mu_2}\right)^2 \text{Var}(\hat{\mu}_2) + \left(\frac{\partial OV_L}{\partial \sigma_1}\right)^2 \text{Var}(\hat{\sigma}_1) + \left(\frac{\partial OV_L}{\partial \sigma_2}\right)^2 \text{Var}(\hat{\sigma}_2),$$

where the derivatives are evaluated using the MLEs

The  $\delta$ -method CIs are

$$\widehat{OV_L} \pm z_{1-\alpha/2} \sqrt{\text{Var}(\widehat{OV_L})}$$

# The Box-Cox transformation

In the binormal framework:

$$X_{ij_i}^{(\lambda)} = \begin{cases} \frac{X_{ij_i}^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \log(X_{ij_i}) & \lambda = 0 \end{cases}$$

The MLE of the common transformation parameter  $\lambda$  can be obtained by maximizing the profile likelihood function given by

$$l(\lambda) = -\frac{n_1}{2} \log \left( \frac{\sum_{i=1}^{n_1} \left( X_{1i}^{(\lambda)} - \frac{\sum_{i=1}^{n_1} X_{1i}^{(\lambda)}}{n_1} \right)^2}{n_1} \right) - \frac{n_2}{2} \log \left( \frac{\sum_{j=1}^{n_2} \left( X_{2j}^{(\lambda)} - \frac{\sum_{j=1}^{n_2} X_{2j}^{(\lambda)}}{n_2} \right)^2}{n_2} \right) + (\lambda - 1) \left( \sum_{i=1}^{n_1} \log X_{1i} + \sum_{j=1}^{n_2} \log X_{2j} \right) + k,$$

where  $k$  is constant.



There are two possible ways to proceed:

- $\delta$ -BC: act as if  $\lambda$  is known, and then use its estimated value to transform the observations and compute the confidence intervals.
- $\delta$ -BC- $\lambda$ : Take into account the variability of the estimated  $\lambda$  by analyzing the full likelihood function. We plug-in  $\hat{\Sigma}^{(\lambda)}$ , an estimator of the covariance matrix of the estimated parameter vector  $(\mu_1^{(\lambda)}, \mu_2^{(\lambda)}, \sigma_1^{(\lambda)}, \sigma_2^{(\lambda)})$ :

$$\text{Var}_{\lambda}(\widehat{OVL}^{(BC)}) \approx \left( \frac{\partial \widehat{OVL}^{(BC)}}{\partial \mu_1}, \frac{\partial \widehat{OVL}^{(BC)}}{\partial \mu_2}, \frac{\partial \widehat{OVL}^{(BC)}}{\partial \sigma_1}, \frac{\partial \widehat{OVL}^{(BC)}}{\partial \sigma_2} \right) \hat{\Sigma}^{(\lambda)} \left( \frac{\partial \widehat{OVL}^{(BC)}}{\partial \mu_1}, \frac{\partial \widehat{OVL}^{(BC)}}{\partial \mu_2}, \frac{\partial \widehat{OVL}^{(BC)}}{\partial \sigma_1}, \frac{\partial \widehat{OVL}^{(BC)}}{\partial \sigma_2} \right)^t$$

# The Logit scale

One can switch to the logit scale,

$$\text{logit}(\widehat{OVL}) \pm z_{1-\alpha/2} \sqrt{\text{Var}(\text{logit}(\widehat{OVL}))},$$

where  $\text{logit}(\widehat{OVL}) = \log(\widehat{OVL} / (1 - \widehat{OVL}))$  and

$$\text{Var}(\text{logit}(\widehat{OVL})) = \text{Var}(\widehat{OVL}) / (\widehat{OVL}(1 - \widehat{OVL}))^2$$

Transforming back to the original scale, confidence intervals for  $OVL$  can be obtained that are properly bounded in  $[0, 1]$ .

We get “L- $\delta$ ”, “L- $\delta$ -BC” and “L- $\delta$ -BC- $\lambda$ ”

# Bootstrap CIs

BC-AN: Another possibility is to consider a parametric bootstrap approach to estimate the variance of the  $OV_L$  estimator:

- 1.- Sample with replacement from  $X_1$  and  $X_2$
- 2.- Box-Cox transformation for each bootstrap sample
- 3.- For  $i = 1, 2$ , calculate  $\hat{\mu}_i$  and  $\hat{\sigma}_i$
- 4.- Calculate  $\widehat{OV_L}$
- 5.- Repeat 1-4  $B$  times.

Then, based on the  $B$  bootstrapped values of  $\widehat{OV_L}$ ,  $\widehat{OV_L}_b^*$ , derive the bootstrap estimate of the variance  $Var_B(\widehat{OV_L})$

$$\widehat{OV_L} \pm z_{1-\alpha/2} \sqrt{Var_B(\widehat{OV_L})}$$

- L-BC-AN: logit scale and back-transformed
- BC-PB: Bootstrap percentile CI
- BC-bias: Bias corrected bootstrap CI

# Non-parametric approaches

The normality assumption may not be satisfied even after applying the Box-Cox transformation. Kernel-based approaches were considered. The density estimator for  $f_{X_1}(x)$  is given by

$$\hat{f}_{X_1}(x) = \frac{1}{n_1} \sum_{i=1}^{n_1} \frac{1}{h} K\left(\frac{x - X_{in_1}}{h}\right)$$

where  $K$  is a kernel function and the bandwidth is given by

$$h = \left(\frac{4}{3}\right)^{1/5} s n_1^{-1/5},$$

with

$$s = \sqrt{\frac{1}{n_1 - 1} \sum_{i=1}^{n_1} \left( X_{in_1} - \sum_{j=1}^{n_1} \frac{X_{jn_1}}{n_1} \right)^2}.$$

Define  $\hat{f}_{X_2}(x)$  analogously and proceed via bootstrap to estimate the variance. We denote this estimator  $\widehat{OVL}^{(K)}$ .

# Basic kernel approach

K-NSR: Confidence intervals for OVL can be obtained via the following steps:

- 1.- Draw bootstrap samples of sizes  $n_1$  and  $n_2$  with replacement from  $X_1$  and  $X_2$ , respectively. Denote these by  $X_{1b}$  and  $X_{2b}$ .
- 2.- Obtain the kernel based OVL estimates, denoted by  $\widehat{OVL}_b^{(K)}$ , for the current ( $b^{th}$ ) bootstrap samples  $X_{1b}$  and  $X_{2b}$ .
- 3.- Repeat Steps 1-2  $B$  times and derive the bootstrap-based estimate of the variance of  $\widehat{OVL}^{(K)}$  by calculating
$$Var\left(\widehat{OVL}^{(K)}\right) = \frac{1}{B-1} \sum_{b=1}^B \left(\widehat{OVL}_b^{(K)} - \overline{\widehat{OVL}}^{(K)}\right)^2$$
where  $\overline{\widehat{OVL}}^{(K)}$  is the mean of the  $B$  terms  $\widehat{OVL}_b^{(K)}$  ( $b = 1, \dots, B$ ).
- 4.- Construct the two-sided  $100(1 - \alpha)\%$  confidence interval of OVL via
$$\overline{\widehat{OVL}}^{(K)} \pm z_{1-\alpha/2} \sqrt{Var\left(\widehat{OVL}^{(K)}\right)}.$$

# Other kernel-based variants

- K-CV: Same but considering the univariate smoothed cross-validation bandwidth selector.
- Working on the logit scale again, the different approaches are “L-K-NSR” or “L-K-CV” depending on the bandwidth.
- Finally, the bootstrap estimates  $\widehat{OVL}_b^*$  can also be used to obtain bootstrap percentile-based confidence intervals (“K-NSR-PB”) and (“K-CV-PB”).

In all these cases we have considered two different kernel functions  $K$ : the Gaussian and the Epanechnikov kernels (and preferred the latter).

# Method overview 1/3

$\delta$	$\widehat{OVL} \pm z_{1-\alpha/2} \sqrt{\widehat{Var}(\widehat{OVL})}$	Parametric approach using the delta method
$\delta$ -BC	$\widehat{OVL}^{(BC)} \pm z_{1-\alpha/2} \sqrt{\widehat{Var}(\widehat{OVL}^{(BC)})}$	Parametric approach using the delta method after the Box-Cox transformation
$\delta$ -BC- $\lambda$	$\widehat{OVL}^{(BC)} \pm z_{1-\alpha/2} \sqrt{\widehat{Var}_{\lambda}(\widehat{OVL}^{(BC)})}$	Parametric approach using the delta method after the Box-Cox transformation taking into account the variability of the estimated transformation parameter
L- $\delta$	$\frac{\exp\left(\logit(\widehat{OVL}) \pm z_{1-\alpha/2} \sqrt{\widehat{Var}(\logit(\widehat{OVL}))}\right)}{1 + \exp\left(\logit(\widehat{OVL}) \pm z_{1-\alpha/2} \sqrt{\widehat{Var}(\logit(\widehat{OVL}))}\right)}$	Parametric approach using the delta method after switching to a logit scale and then transforming back
L- $\delta$ -BC	$\frac{\exp\left(\logit(\widehat{OVL}^{(BC)}) \pm z_{1-\alpha/2} \sqrt{\widehat{Var}(\logit(\widehat{OVL}^{(BC)}))}\right)}{1 + \exp\left(\logit(\widehat{OVL}^{(BC)}) \pm z_{1-\alpha/2} \sqrt{\widehat{Var}(\logit(\widehat{OVL}^{(BC)}))}\right)}$	Parametric approach using the delta method after the Box-Cox transformation after switching to a logit scale and then transforming back

# Method overview 2/3

L-δ-BC-λ	$\frac{\exp\left(\text{logit}\left(\widehat{OVL}^{(BC)}\right) \pm z_{1-\alpha/2} \sqrt{\text{Var}_\lambda\left(\text{logit}\left(\widehat{OVL}^{(BC)}\right)\right)}\right)}{1 + \exp\left(\text{logit}\left(\widehat{OVL}^{(BC)}\right) \pm z_{1-\alpha/2} \sqrt{\text{Var}_\lambda\left(\text{logit}\left(\widehat{OVL}^{(BC)}\right)\right)}\right)}$	<p>Parametric approach using the delta method after the</p> <p>Box-Cox transformation in the logit scale and back-transformed</p> <p>considering the variability of the estimated transformation parameter</p>
BC-AN	$\widehat{OVL}^{(BC)} \pm z_{1-\alpha/2} \sqrt{\text{Var}_B\left(\widehat{OVL}^{(BC)}\right)}$	<p>Parametric approach</p> <p>using a bootstrap-based approach to estimate the variance</p>
BC-PB	$\left(\widehat{OVL}_{(1-\alpha/2)}^*, \widehat{OVL}_{(\alpha/2)}^*\right)$	<p>Parametric approach using a bootstrap percentile approach</p>
BC-bias	$\left(\widehat{OVL}_{(\alpha_1)}^*, \widehat{OVL}_{(\alpha_2)}^*\right)$	<p>Parametric approach using a bootstrap bias-corrected approach</p>

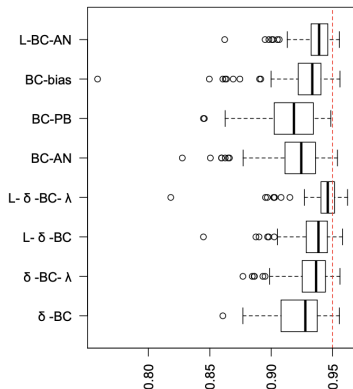


# Method overview 3/3

L-BC-AN	$\frac{\exp\left(\text{logit}\left(\widehat{OVL}^{(BC)}\right) \pm z_{1-\alpha/2} \sqrt{\text{Var}_B\left(\text{logit}\left(\widehat{OVL}^{(BC)}\right)\right)}\right)}{1 + \exp\left(\text{logit}\left(\widehat{OVL}^{(BC)}\right) \pm z_{1-\alpha/2} \sqrt{\text{Var}_B\left(\text{logit}\left(\widehat{OVL}^{(BC)}\right)\right)}\right)}$	BC-AN procedure carried out in the logit scale and back-transformed
K-NSR K-CV	$\widehat{OVL}^{(K)} \pm z_{1-\alpha/2} \sqrt{\text{Var}\left(\widehat{OVL}^{(K)}\right)}$	Kernel approach estimating the variance via bootstrap
L-K-NSR L-K-CV	$\frac{\exp\left(\text{logit}\left(\widehat{OVL}^{(K)}\right) \pm z_{1-\alpha/2} \sqrt{\text{Var}\left(\text{logit}\left(\widehat{OVL}^{(K)}\right)\right)}\right)}{1 + \exp\left(\text{logit}\left(\widehat{OVL}^{(K)}\right) \pm z_{1-\alpha/2} \sqrt{\text{Var}\left(\text{logit}\left(\widehat{OVL}^{(K)}\right)\right)}\right)}$	Kernel approach estimating the variance via bootstrap in the logit scale and back-transformed
K-NSR-PB K-CV-PB	$\left(\widehat{OVL}_{(\alpha_1)}^{(K)*}, \widehat{OVL}_{(\alpha_2)}^{(K)*}\right)$	Kernel approach using a bootstrap percentile approach

# Some simulations results

All parametric simulations considered...



# R-package OVL.CI

## OVL.CI: Inference on the Overlap Coefficient: The Binormal Approach and Alternatives

Provides functions to construct confidence intervals for the Overlap Coefficient (OVL). OVL measures the similarity between two distributions through the overlapping area of their distribution functions. Given its intuitive description and ease of visual representation by the straightforward depiction of the amount of overlap between the two corresponding histograms based on samples of measurements from each one of the two distributions, the development of accurate methods for confidence interval construction can be useful for applied researchers. Implements methods based on the work of Franco-Pereira, A.M., Nakas, C.T., Reiser, B., and Pardo, M.C. (2021) <[doi:10.1177/09622802211046386](https://doi.org/10.1177/09622802211046386)>.

Version: 0.1.0  
Depends: R (≥ 2.10)  
Imports: [ks](#)  
Suggests: [testthat](#) (≥ 3.0.0)  
Published: 2023-11-13  
DOI: [10.32614/CRAN.package.OVL.CI](https://doi.org/10.32614/CRAN.package.OVL.CI)  
Author: Alba M. Franco-Pereira [aut, cre, cph], Christos T. Nakas [aut], Benjamin Reiser [aut], M.Carmen Pardo [aut]  
Maintainer: Alba M. Franco-Pereira <[albfranc@ucm.es](mailto:albfranc@ucm.es)>  
License: [GPL-2](#)  
NeedsCompilation: no  
Language: en-US  
Materials: [NEWS](#)  
CRAN checks: [OVL.CI results](#)

Documentation:

Reference manual: [OVL.CI.pdf](#)

Downloads:

Package source: [OVL.CI\\_0.1.0.tar.gz](#)

Windows binaries: r-devel: [OVL.CI\\_0.1.0.zip](#), r-release: [OVL.CI\\_0.1.0.zip](#), r-oldrel: [OVL.CI\\_0.1.0.zip](#)

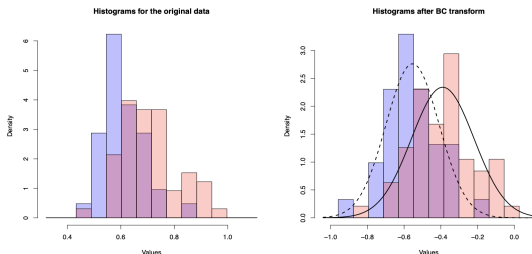
macOS binaries: r-release (arm64): [OVL.CI\\_0.1.0.tgz](#), r-oldrel (arm64): [OVL.CI\\_0.1.0.tgz](#), r-release (x86\_64): [OVL.CI\\_0.1.0.tgz](#), r-oldrel (x86\_64): [OVL.CI\\_0.1.0.tgz](#)

Linking:

Please use the canonical form <https://CRAN.R-project.org/package=OVL.CI> to link to this page.

# Illustration


Myoinositol (MI) over creatine (Cr) measured via proton MRS was used from a published study involving 39 neurologically asymptomatic (NAS) HIV+ patients and 37 HIV- controls (NEG), as a biomarker for dementia.



Method	Point Estimate	95% CI		Width
		LC	UC	
L- $\delta$ -BC- $\lambda$	0.5751	0.4390	0.7313	0.2923
L-BC-AN	0.5934	0.4360	0.7337	0.2977

# ROC analysis?

Nakas, C.T., Bantis, L.E., & Gatsonis, C.A. (2023). ROC Analysis for Classification and Prediction in Practice (1st ed.). Chapman and Hall/CRC. <https://doi.org/10.1201/9780429170140>


 Taylor & Francis Group  
an informa business

T&F eBooks ▾

Search for keywords, authors, titles, ISBN

About Us ▾ Subjects ▾ Browse ▾ Products ▾ Request a trial Librarian Resou

Home > Mathematics & Statistics > Statistics & Probability > Statistics > Statistical Theory & Methods > ROC




Book

## ROC Analysis for Classification and Prediction in Practice

By *Christos T Nakas, Leonidas E Bantis, Constantine A Gatsonis*

Edition	1st Edition
First Published	2023
eBook Published	15 May 2023
Pub. Location	New York
Imprint	Chapman and Hall/CRC
DOI	<a href="https://doi.org/10.1201/9780429170140">https://doi.org/10.1201/9780429170140</a>
Pages	234
eBook ISBN	9780429170140
Subjects	Mathematics & Statistics, Medicine, Dentistry, Nursing & Allied Health

 Accessibility Information

# Why Covariate Adjustment?

- OVL is often affected by patient characteristics (e.g. age) since covariates may significantly influence distributional overlap.
- Ignoring covariates can mask heterogeneity.
- Need to estimate  $OVL(z)$ : overlap conditional on covariates using linear regression with a possible Box-Cox transformation.
- We develop procedures for CI calculation for  $OVL(z)$ .

Accepted in *Journal of Biopharmaceutical Statistics*

## Confidence intervals for the covariate-specific overlap coefficient (OVL)

M. Carmen Pardo<sup>1,2</sup>, Alba M. Franco-Pereira<sup>1,2,\*</sup>, Benjamin Reiser<sup>3</sup> and Christos T. Nakas<sup>4,5</sup>

# Model Framework

- Linear regression framework:  $X_i = \beta^T Z + \varepsilon$
- Conditional densities  $f_1(x|z), f_2(x|z)$
- Box-Cox transformation to normalize data if needed
- $\text{OVL}(z)$ : overlap between conditional densities



# Basic model components

- Let  $\left\{ \left( z_{\overline{D}_i}, x_{\overline{D}_i} \right) \right\}_{i=1}^{n_{\overline{D}}}$  and  $\left\{ \left( z_{D_i}, x_{D_i} \right) \right\}_{i=1}^{n_D}$  be two independent random samples of test outcomes and covariates from non-diseased and diseased groups of size  $n_{\overline{D}}$  and  $n_D$ , respectively.
- Where  $z_{\overline{D}_i} = \left( z_{\overline{D}_{i,1}}, \dots, z_{\overline{D}_{i,p_1}} \right)^t$  and  $z_{D_i} = \left( z_{D_{i,1}}, \dots, z_{D_{i,p_2}} \right)^t$  be  $p_1$ -dimensional and  $p_2$ -dimensional vectors of covariates, respectively.
- Assume that the (possibly transformed) non-diseased and diseased marker values can be described as linear functions of the explanatory variables.

$$X_{\overline{D}} = Z_{\overline{D}}\beta_{\overline{D}} + \epsilon_{\overline{D}},$$

$$X_D = Z_D\beta_D + \epsilon_D,$$

- where  $X_{\overline{D}} = \left( x_{\overline{D}_1}, \dots, x_{\overline{D}_{n_{\overline{D}}}} \right)^t$ ,  $X_D = \left( x_{D_1}, \dots, x_{D_{n_D}} \right)^t$  and  $\beta_{\overline{D}}$  and  $\beta_D$  are the parameters' column vectors of size  $p_1$  and  $p_2$ , respectively.

# OVL(z) estimation 1/2

The covariate adjusted *OVL* for a vector of covariates  $Z = z$  of dimension  $p$  is given by

$$OVL(z) = 1 + \Phi\left(\frac{x_1(z) - \beta_D^t z}{\sigma_{\bar{D}}}\right) - \Phi\left(\frac{x_1(z) - \beta_D^t z}{\sigma_D}\right) - \Phi\left(\frac{x_2(z) - \beta_D^t z}{\sigma_{\bar{D}}}\right) + \Phi\left(\frac{x_2(z) - \beta_D^t z}{\sigma_D}\right)$$

assuming that  $\sigma_{\bar{D}}^2 < \sigma_D^2$ , where:

$$x_1(z) = \frac{(\beta_D^t z \sigma_D^2 - \beta_D^t z \sigma_{\bar{D}}^2) - \sigma_{\bar{D}} \sigma_D \sqrt{(\beta_D^t z - \beta_D^t z)^2 + (\sigma_{\bar{D}}^2 - \sigma_D^2) \log\left(\frac{\sigma_{\bar{D}}^2}{\sigma_D^2}\right)}}{(\sigma_D^2 - \sigma_{\bar{D}}^2)}$$

and

$$x_2(z) = \frac{(\beta_D^t z \sigma_D^2 - \beta_D^t z \sigma_{\bar{D}}^2) + \sigma_{\bar{D}} \sigma_D \sqrt{(\beta_D^t z - \beta_D^t z)^2 + (\sigma_{\bar{D}}^2 - \sigma_D^2) \log\left(\frac{\sigma_{\bar{D}}^2}{\sigma_D^2}\right)}}{(\sigma_D^2 - \sigma_{\bar{D}}^2)}.$$

## OVL(z) estimation 2/2

$OVL(z)$  is estimated by plugging-in the relevant estimates for the unknown  $\beta_{\bar{D}}, \beta_D, \sigma_{\bar{D}}^2, \sigma_D^2$  in the formulae given in the previous slide:

$$\hat{\beta}_{\bar{D}} = \left( Z_{\bar{D}}^t Z_{\bar{D}} \right)^{-1} Z_{\bar{D}}^t X_{\bar{D}}; \quad \hat{\beta}_D = \left( Z_D^t Z_D \right)^{-1} Z_D^t X_D;$$

$$\hat{\sigma}_{\bar{D}}^2 = \frac{\left( X_{\bar{D}}^t X_{\bar{D}} - \hat{\beta}_{\bar{D}}^t Z_{\bar{D}}^t X_{\bar{D}} \right)}{n_{\bar{D}} - p}; \quad \hat{\sigma}_D^2 = \frac{\left( X_D^t X_D - \hat{\beta}_D^t Z_D^t X_D \right)}{n_D - p}.$$

# Confidence Intervals for $OV_L(z)$

- Delta method: analytic variance (complex)
- Bootstrap methods:
  - PB (percentile)
  - Bias-corrected (bias)
  - Logit + AN using bootstrap variance estimate (L-AN, LD-AN)
  - Box-Cox extensions (L-BC-AN, LD-BC-AN, BC-PB, BC-bias)
- Logit transformation helps bound CI in  $[0,1]$

# Basic bootstrap CIs for OVL(z)

Method L-AN:

- Compute  $\hat{\beta}_{\overline{D}}^t, \hat{\beta}_D^t, \hat{\sigma}_{\overline{D}}^2, \hat{\sigma}_D^2$  from the original sample to calculate  $\widehat{OVL}(z)$
- Sample with replacement from  $(X_{\overline{D}}, Z_{\overline{D}})$  and  $(X_D, Z_D)$  and calculate  $\widehat{OVL}(z)$  for the bootstrap sample.
- Repeat  $B$  times. Then, based on the  $B$  bootstrapped values  $\widehat{OVL}_b(z)$ , derive the bootstrap estimate of  $Var_B(\widehat{OVL}(z))$ .
- Construct the two-sided  $100(1 - \alpha)\%$  confidence interval of  $OVL(z)$  switching to the logit scale and then transforming back to the original scale, through

$$\text{logit}(\widehat{OVL}(z)) \pm z_{1-\alpha/2} \sqrt{\text{Var}(\text{logit}(\widehat{OVL}(z)))}.$$

## Alternatives:

- LD-AN: Estimate the variance of the logit of  $OVL$  using the bootstrap estimate of logit  $OVL$  directly.
- The bootstrap estimates,  $\widehat{OVL}_b(z)$  can also be used to obtain bootstrap percentile confidence intervals ('PB') as well as the bias corrected bootstrap confidence interval ('bias').

Box-Cox transformation:

$$x_{\overline{D}}^{(\lambda)} | z = \begin{cases} \frac{(x_{\overline{D}}^{(\lambda)} | z) - 1}{\lambda} & \lambda \neq 0 \\ \log(x_{\overline{D}}^{(\lambda)} | z) & \lambda = 0 \end{cases} \quad x_D^{(\lambda)} | z = \begin{cases} \frac{(x_D^{(\lambda)} | z) - 1}{\lambda} & \lambda \neq 0 \\ \log(x_D^{(\lambda)} | z) & \lambda = 0 \end{cases}$$

where the maximum likelihood estimate of the common transformation parameter  $\lambda$  can be obtained by maximizing the log-profile likelihood function given by

$$l(\lambda) = \sum_{i \in \{\overline{D}, D\}} \left[ -\frac{n_i}{2} (\log 2\pi \hat{\sigma}_i^2) - \frac{1}{2} \sum_{j=1}^{n_i} \left( (x_{ij}^{(\lambda)} - \hat{\beta}_i^t z_i)^2 / \hat{\sigma}_i^2 \right) + (\lambda - 1) \sum_{j=1}^{n_i} \log x_{ij}^{(\lambda)} \right]$$

- When the Box-Cox transformation is used for the regression model, we add a step 0 in the bootstrap algorithm using the transformation first on the original sample data.
- To obtain the bootstrapped estimates and compute the variance, we need to carry out a Box-Cox transformation between steps 2 and 3. Specifically, after sampling with replacement from the original sample, we apply the BC transformation.
- Methods denoted by L-BC-AN, LD-BC-AN, BC-PB and BC-bias then arise.
- The Box-Cox transformation leads to a nonlinear relationship between the marker of interest and the covariates.

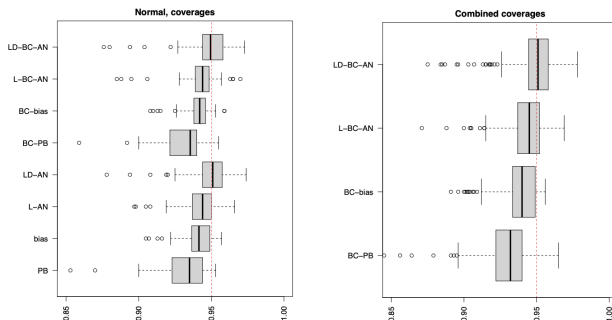


# Simulation Study Design

- Scenarios: Normal, LogNormal, PowerNormal
- Sample sizes: 50, 100, 200, 500
- Metrics: Coverage probability, CI width
- Methods compared: PB, bias, L-AN, LD-AN, L-BC-AN, LD-BC-AN, BC-PB, BC-bias

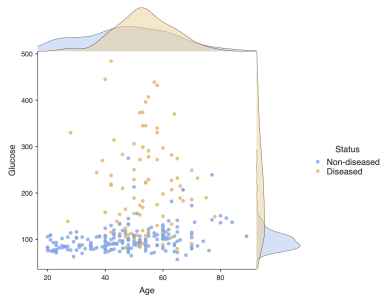
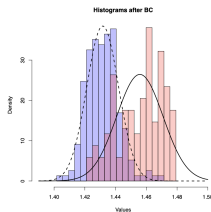
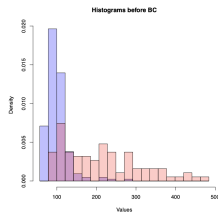
# Key Result: LD-BC-AN Performs Best overall

- LD-BC-AN: stable, accurate, and robust
- Box-Cox adds robustness with minimal cost under normality
- PB often undercovers; logit-based methods superior



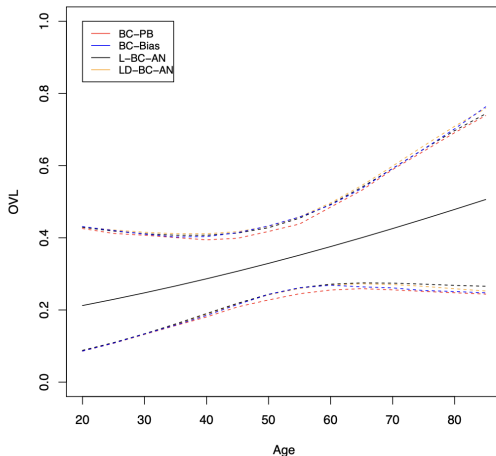
# Real Data: Glucose as a Diabetes marker adjusted for Age

- 286 subjects (198 non-diseased, 88 diseased)
- Post-prandial glucose as marker for diabetes
- Age as covariate



# OVL(Age): Covariate-Specific Insight

- Better discrimination in younger subjects
- Age-adjusted OVL provides meaningful stratified insight



- Introduced inference for covariate-specific OVL
- Logit + bootstrap-based methods work best
- Box-Cox transformation improves flexibility
- LD-BC-AN is a recommended method
- Future steps: Incorporate into R package ('OVL.CI')

# Thank You

Questions?

`cnakas@uth.gr`