

Second-Generation p -Values in a High Dimensional Analysis of Prostate Cancer Variants

Valerie Welty
Department of Biostatistics
Vanderbilt University



THE AMERICAN STATISTICIAN
2019, VOL. 73, NO. 51, 157–167: Statistical Inference in the 21st Century
<https://doi.org/10.1080/00031305.2019.1537893>



 OPEN ACCESS



An Introduction to Second-Generation p -Values

Jeffrey D. Blume^{a,b}, Robert A. Greevy^a, Valerie F. Welty^a, Jeffrey R. Smith^{c,d}, and William D. Dupont^{a,e}

^aDepartment of Biostatistics, Vanderbilt University School of Medicine, Nashville, TN; ^bDepartment of Biomedical Informatics, Vanderbilt University School of Medicine, Nashville, TN; ^cDepartment of Medicine, Vanderbilt University School of Medicine, Nashville, TN; ^dDepartment of Cancer Biology, Vanderbilt University School of Medicine, Nashville, TN; ^eDepartment of Health Policy, Vanderbilt University School of Medicine, Nashville, TN

ABSTRACT

Second generation p -values preserve the simplicity that has made p -values popular while resolving critical flaws that promote misinterpretation of data, distraction by trivial effects, and unreproducible assessments of data. The second-generation p -value (SGPV) is an extension that formally accounts for scientific relevance by using a composite null hypothesis that captures null and scientifically trivial effects. Because the majority of spurious findings are small effects that are technically nonnull but practically indistinguishable from the null, the second-generation approach greatly reduces the likelihood of a false discovery. SGPVs promote transparency, rigor and reproducibility of scientific results by a priori identifying which candidate hypotheses are practically meaningful and by providing a more reliable statistical summary of when the data are compatible with the candidate hypotheses or null hypotheses, or when the data are inconclusive. We illustrate the importance of these advances using a dataset of 247,000 single-nucleotide polymorphisms, i.e., genetic markers that are potentially associated with prostate cancer.

ARTICLE HISTORY

Received March 2018
Revised September 2018

KEYWORDS

Likelihood ratios; Null hypothesis; p -Value; Statistical evidence

Typical concerns with standard p -value approaches

- Statistical significance \neq clinical or practical significance
- Large p -values do not indicate support for the null hypothesis
- p -value adjustments for multiple comparisons
 - Often conservative
 - No universal solution
- Ranking findings by p -value may miss interesting large effects

Alternatives

- Consider an interval null hypothesis
- Differentiate between when the data support:
 - Only alternative hypotheses
 - Only null hypotheses
 - Null and alternative hypotheses (inconclusive)
- Type I Error rate shrinks as the sample size increases
- Rank findings by clinical and statistical importance

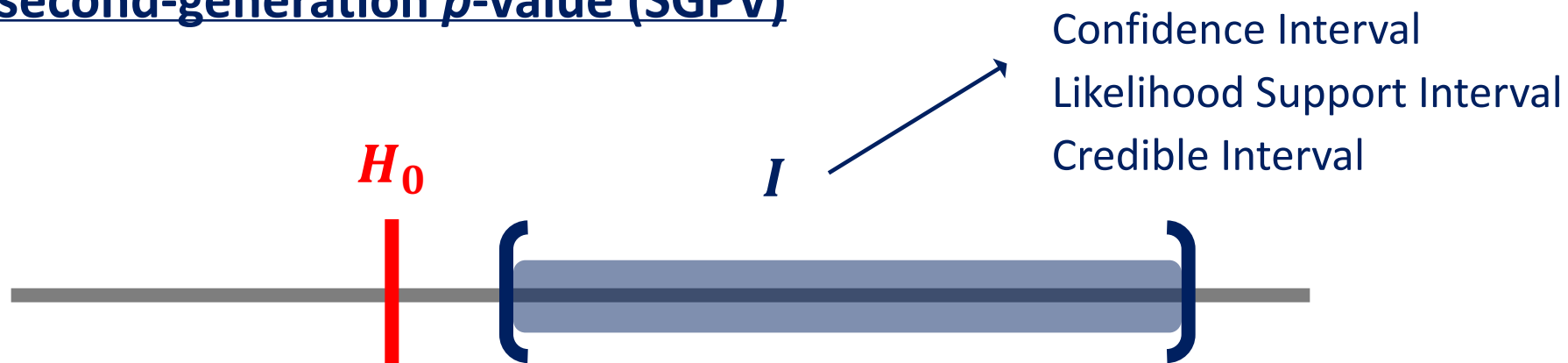
Alternatives

- Consider an interval null hypothesis
- Differentiate between when the data support:
 - Only alternative hypotheses
 - Only null hypotheses
 - Null and alternative hypotheses (inconclusive)
- Type I Error rate shrinks as the sample size increases
- Rank findings by clinical and statistical importance

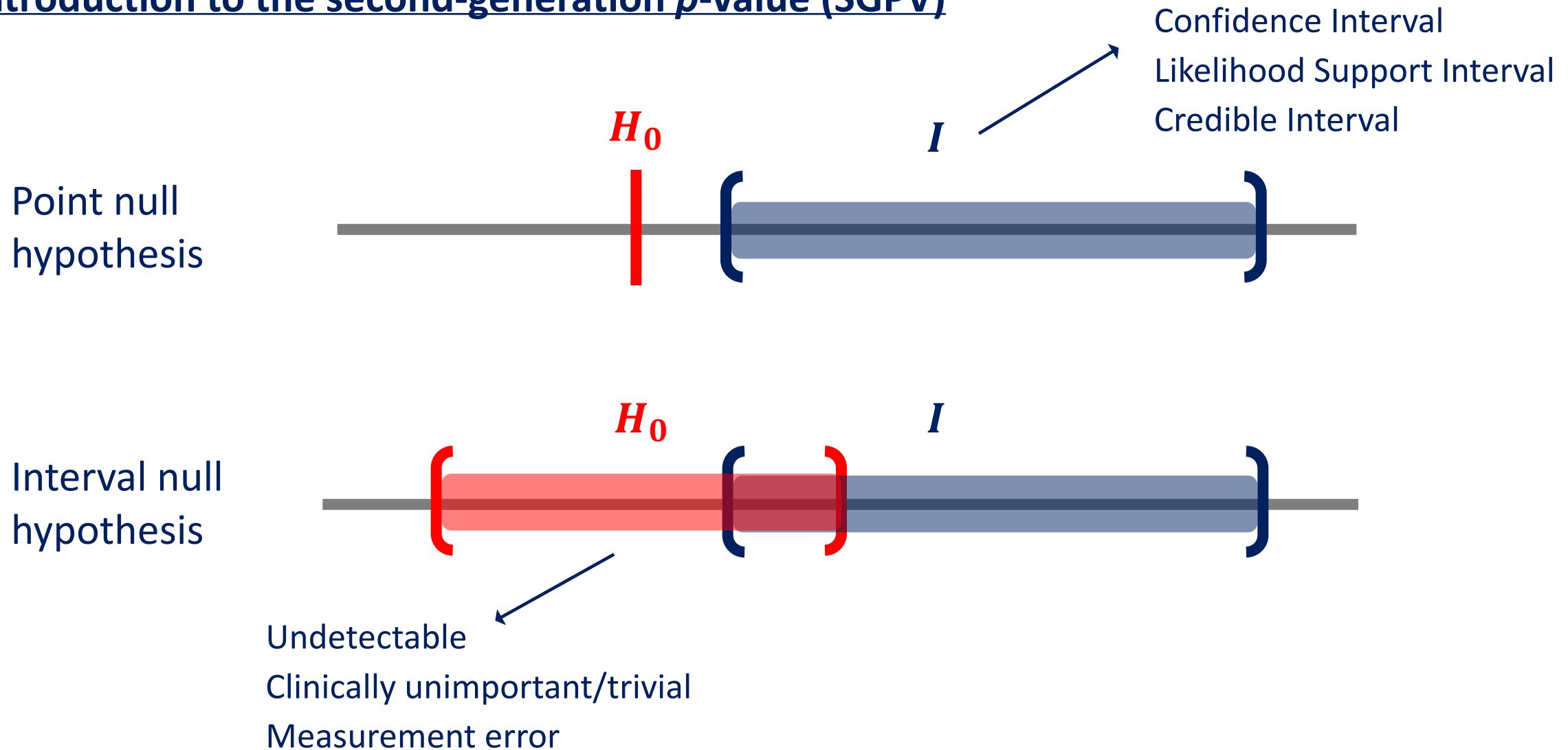
second-generation
 p -value

Introduction to the second-generation p -value (SGPV)

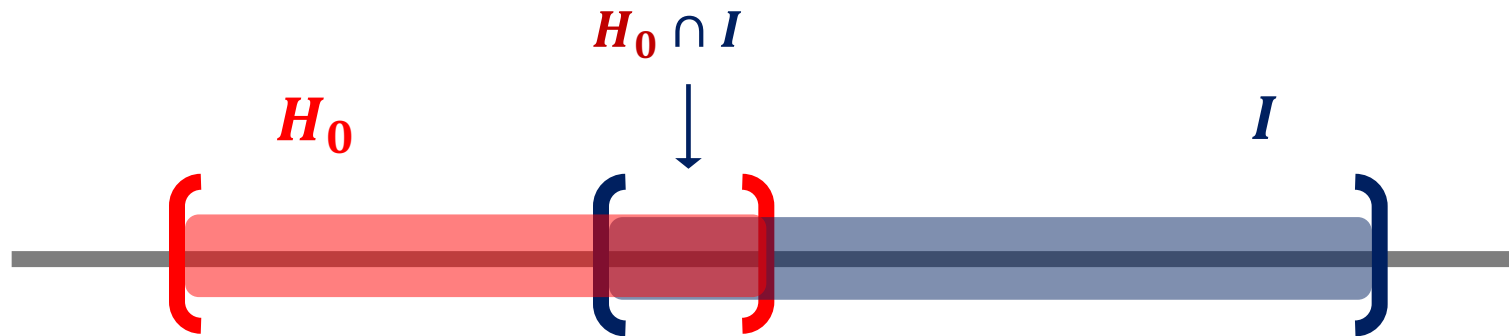
Point null
hypothesis



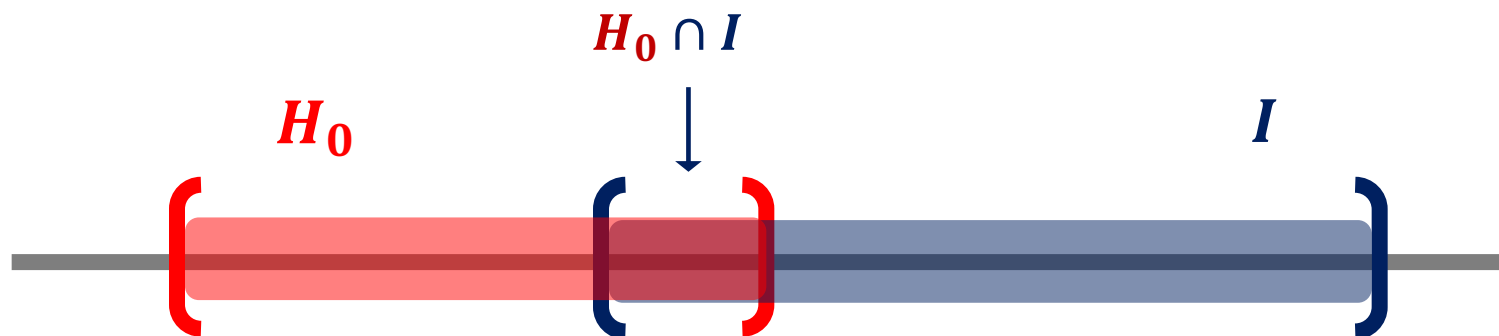
Introduction to the second-generation p -value (SGPV)



Introduction to the second-generation p -value (SGPV)



Introduction to the second-generation p -value (SGPV)



$$p_{\delta} = \frac{|I \cap H_0|}{|I|} \times \max \left\{ \frac{|I|}{2|H_0|}, 1 \right\}$$

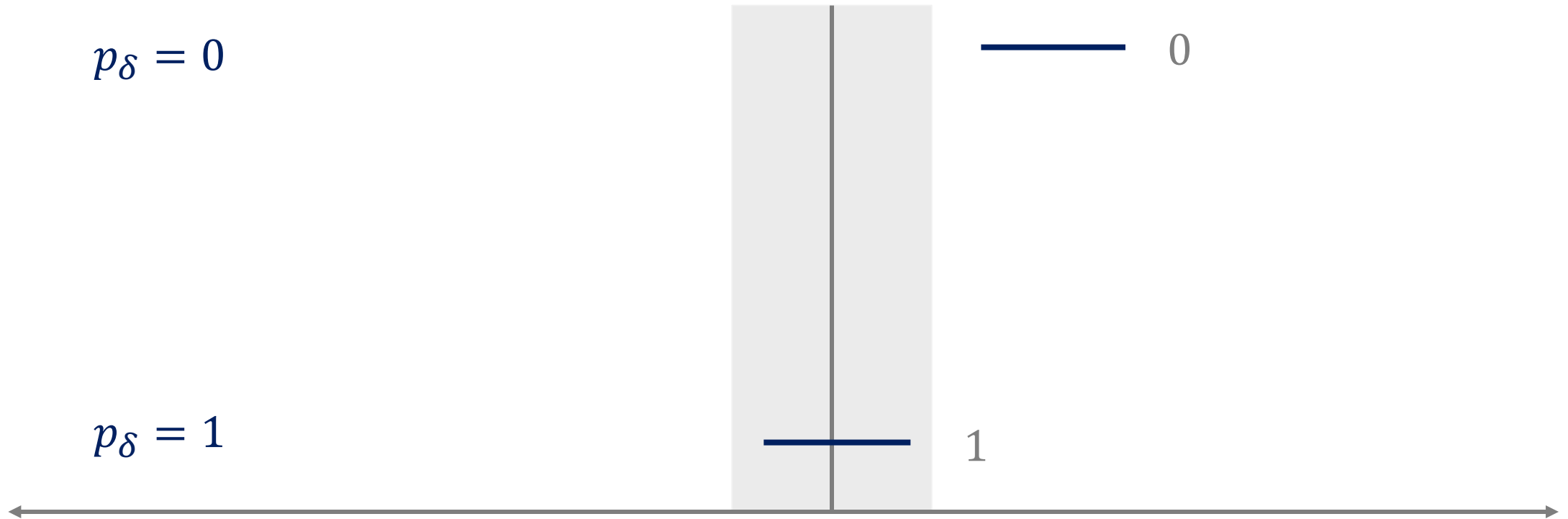
Proportion of data-supported hypotheses that are also null hypotheses

Small-sample correction factor

shrinks proportion towards $\frac{1}{2}$ when $|I|$ is greater than $2|H_0|$

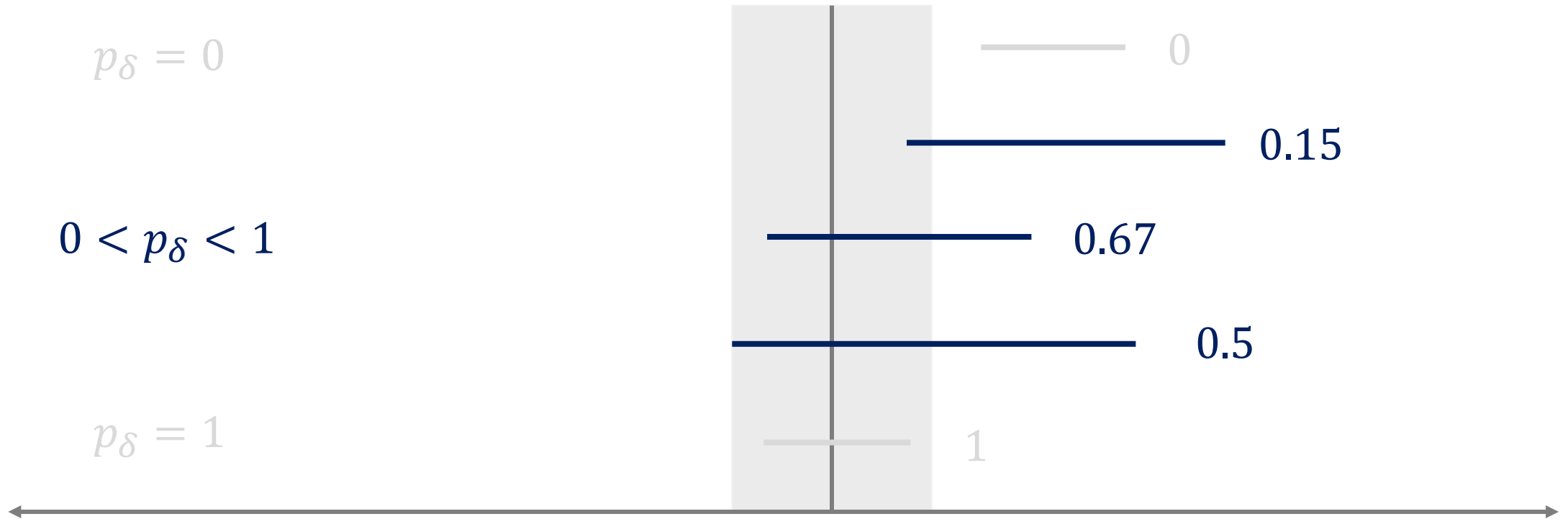
Introduction to the second-generation p -value (SGPV)

- Provides a single-number summary of when the data support only null effects, only meaningful alternative effects, or are inconclusive



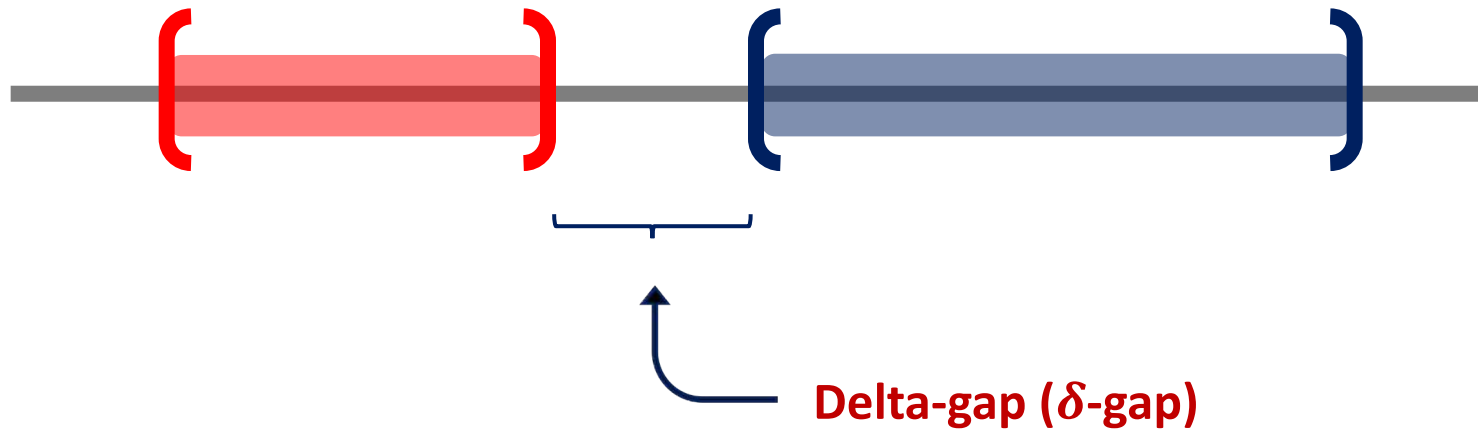
Introduction to the second-generation p -value (SGPV)

- Provides a single-number summary of when the data support only null effects, only meaningful alternative effects, or are inconclusive



Ranking second-generation p -values

- When $p_\delta = 0$, there is a gap between the intervals. The length of that gap is the **delta-gap**



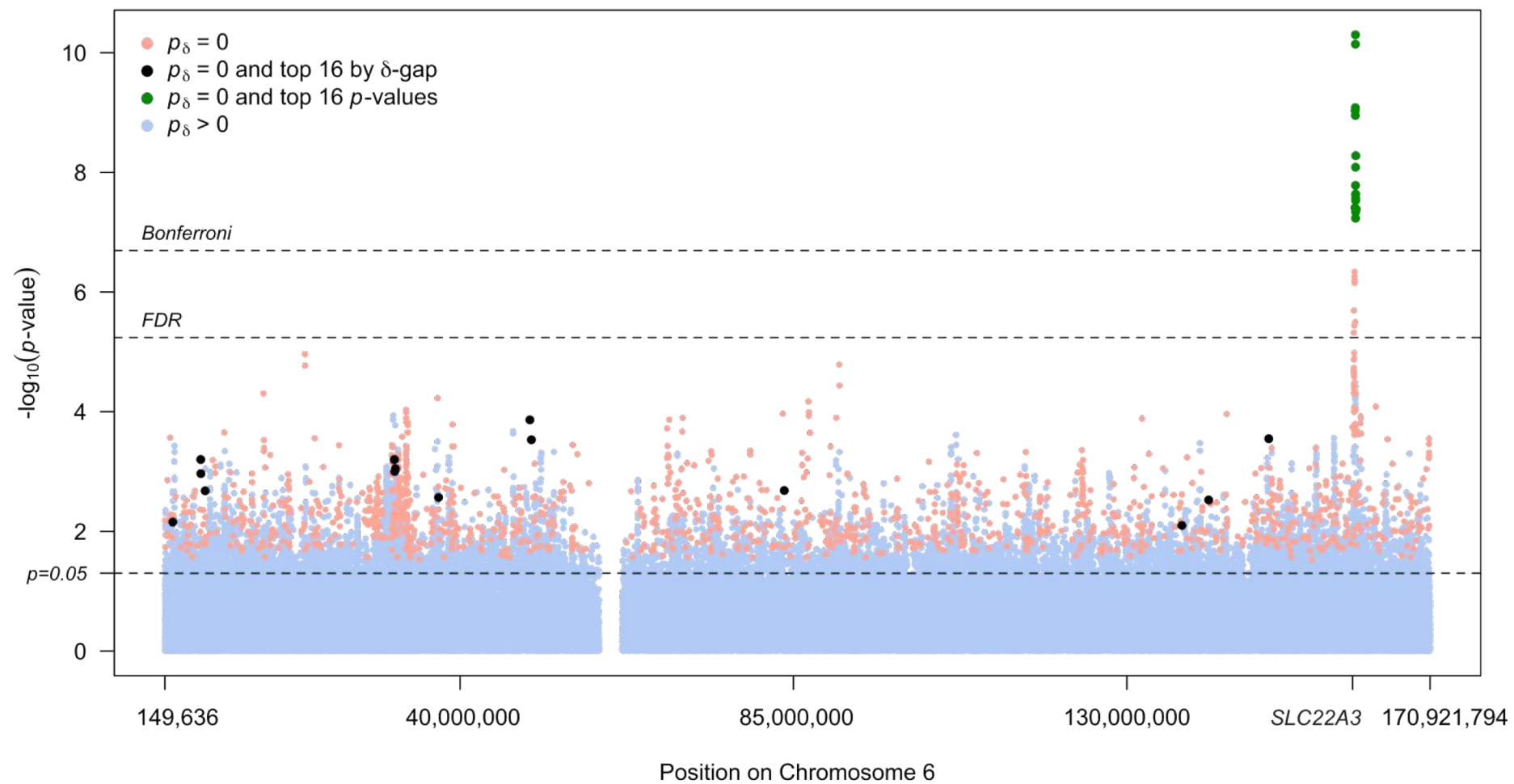
- Allows refined differentiation among findings

Example with prostate cancer variants

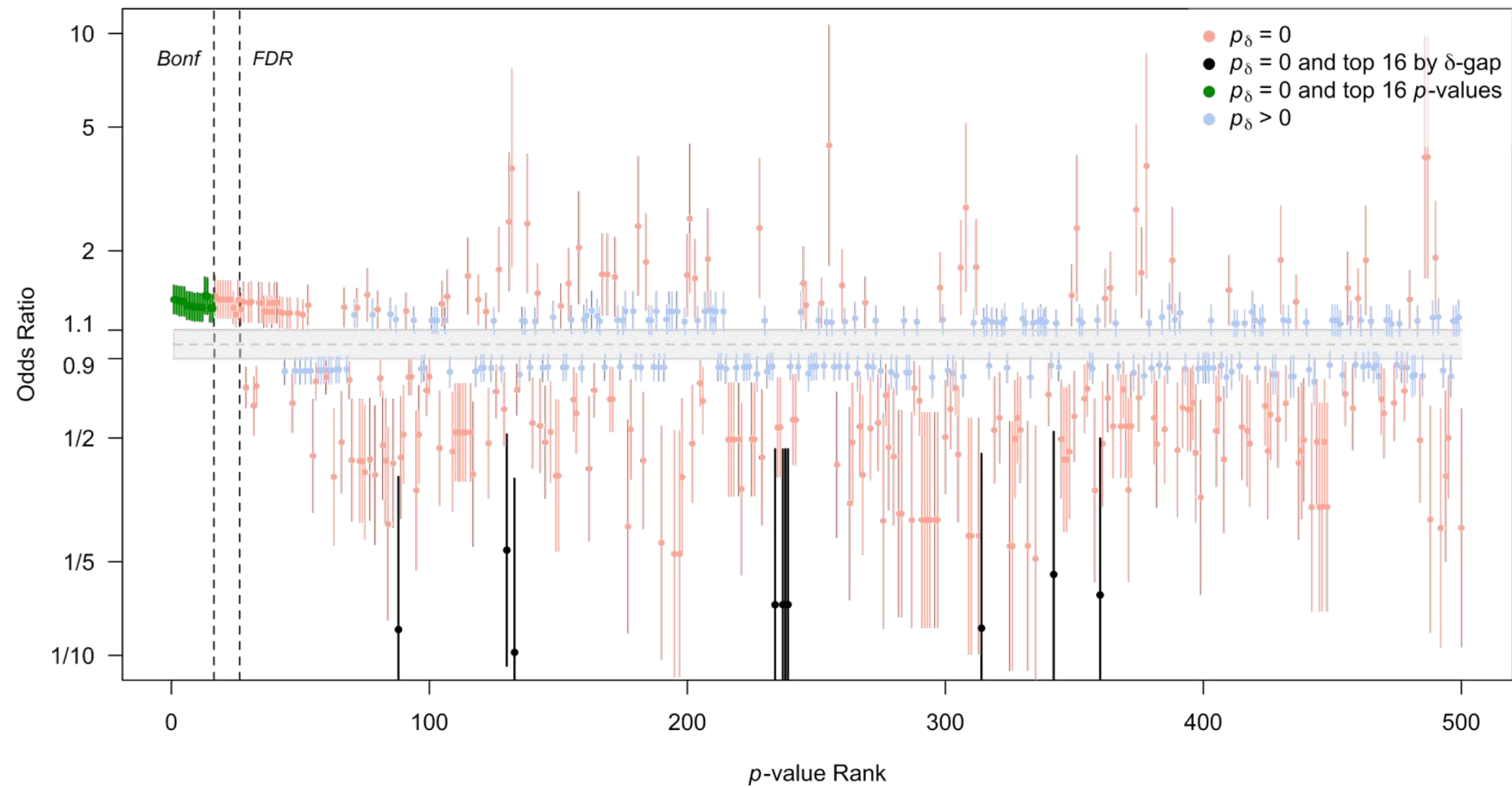
- International Consortium for Prostate Cancer Genetics (Schaid and Chang 2015; ICPCG 2018)
 - ~ 3,900 subjects (2,500 cases & 1,400 controls)
 - ~ 247,000 single-nucleotide polymorphisms (SNPs) from Chromosome 6
- Goal: Identify ‘interesting’ SNPs potentially associated with prostate cancer

$$H_0: OR \in [0.9, 1.11]$$

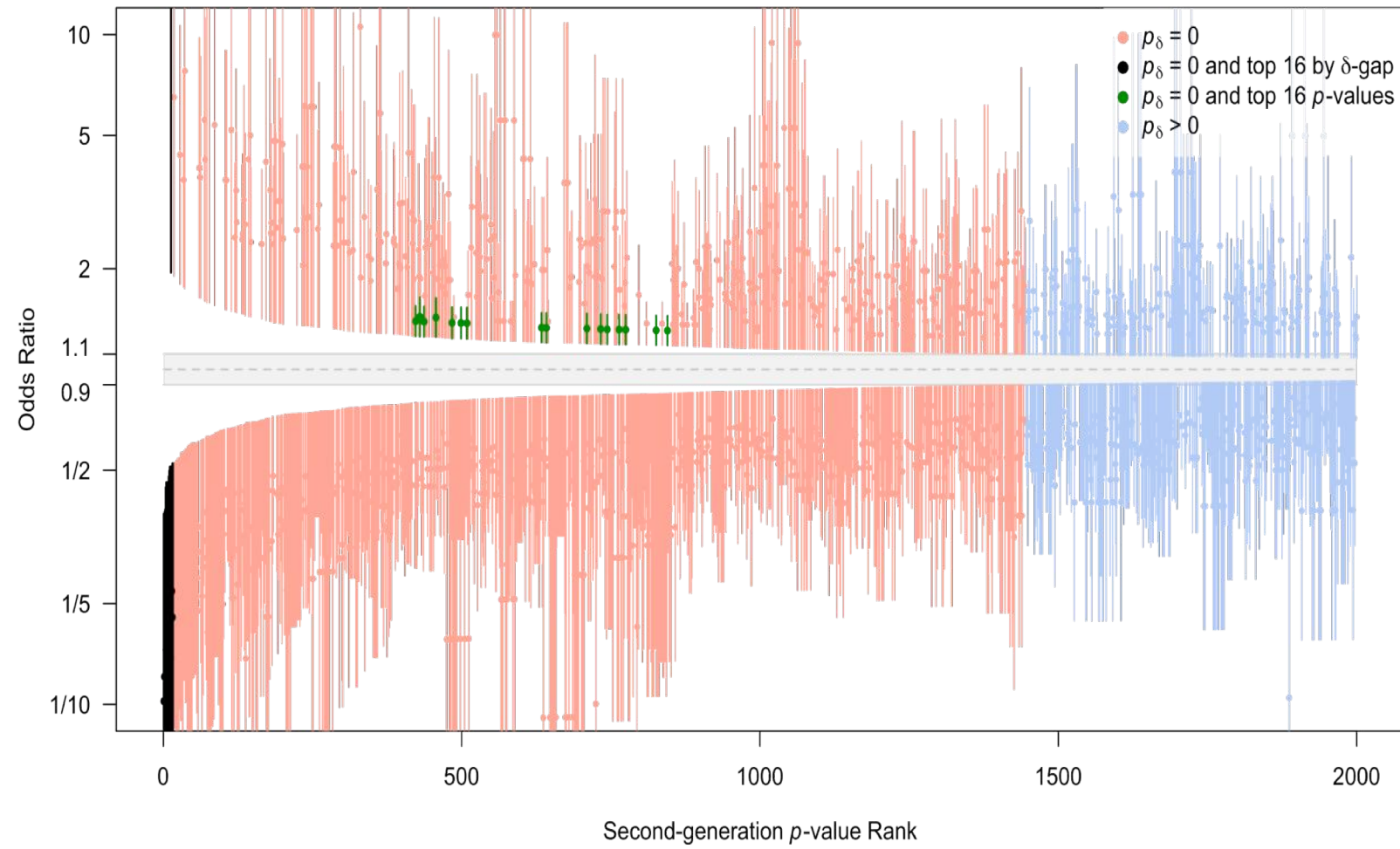
Manhattan Plot



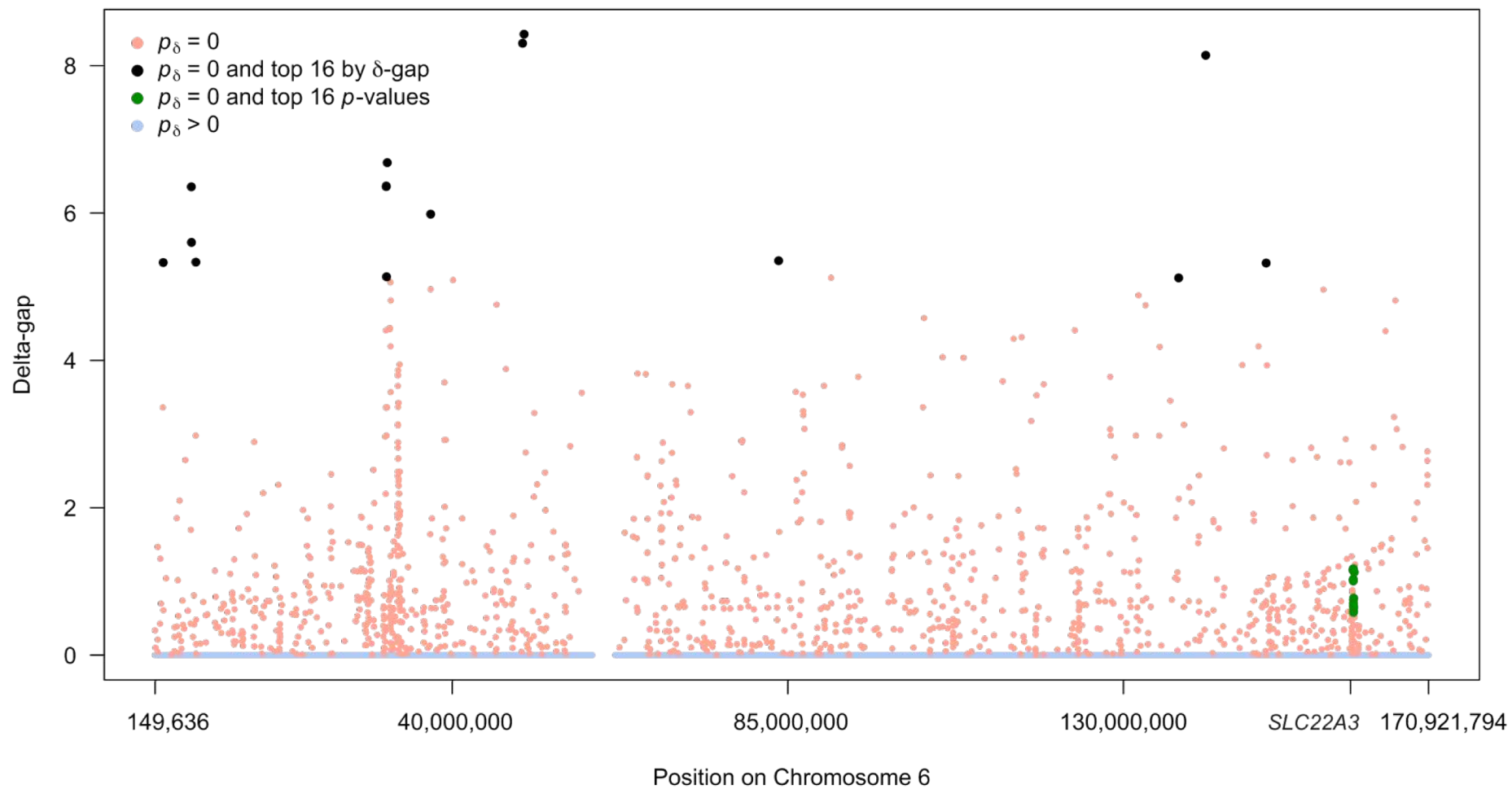
Top 500 support intervals by p -value ranking



Top 2,000 support intervals by SGPV ranking



SGPV Manhattan plot




Discussed in papers (Blume et al. 2018, Blume et al. 2019)

- SGPV error rate profile
 - Both Type I and II Error rates converge to 0 with sample size
- False discovery rate profile
 - Maintains FDR improvement associated with multiple comparisons adjustments while lowering the false confirmation rate

Ongoing work

- SGPV false discovery rates (FDR_δ)
 - Estimation methods
 - Rank SGPV findings by FDR_δ

Recommendations

- Report all finding where $p_\delta = 0$ and $p_\delta = 1$ (“discoveries” and “confirmations”)
- Identify top findings among those with $p_\delta = 0$ by large delta-gaps or low FDR_δ values
 - Prioritizes clinical impact over precision
 - Incorporates reliability of findings
- Above set of findings is very different from those of classical p -value approaches
- R package available at  github.com/weltybiostat/sgpv

`sgpvalue`
`sgpower`

`plotsgpv`
`fdrisk`

Support and References

TREAT Research Group (Dept. of Thoracic Surgery) and Statistical Evidence in Data Science (SEDS) Lab

- Jeffrey Blume (PI)  statisticalevi.dence.com
- Thomas Stewart
- Megan Hollister

Blume JD, Greevy RA, Welty VF, Smith JR, Dupont WD (2019) An Introduction to Second-Generation p -Values, The American Statistician, 73:sup1, 157-167, DOI: [10.1080/00031305.2018.1537893](https://doi.org/10.1080/00031305.2018.1537893)

Blume JD, D'Agostino McGowan L, Dupont WD, Greevy RA Jr (2018) Second-generation p -values: Improved rigor, reproducibility, & transparency in statistical analyses. PLoS ONE 13(3): e0188299. <https://doi.org/10.1371/journal.pone.0188299>