1    **Supporting information for: Unbiased estimation of the number of segregating sites across**

2    **unequal sample sizes**

3    William Hemstrom*, Mark R. Christie*

4

5    **Corresponding author information:**
6    * William Hemstrom; email: whemstro@purdue.edu
7    * Mark Christie; email: christ99@purdue.edu

8
9    **This PDF includes:**

12 **Supplementary Example:**

13 The unbiased expected number of segregating sites after standardizing sample sizes, $E(N_S)$, can

14 be calculated with $\gamma$ automatically set to the smallest sample size across all sample groups

15 automatically across multiple sample groups using the R package "snpR" (W. Hemstrom &

16 Jones, 2023). First, install snpR if needed (this can be skipped if already installed) and load it:

```r
# Install and load snpR. Comment out first two lines if not needed.

install.packages("remotes")
remotes::install_github("hemstrow/snpR")
library(snpR)
```

22 Next, download the example data:

```r
# download the metadata (if copy/pasting, check line breaks)
meta <- read.table(url("https://raw.githubusercontent.com/ChristieLab/seg_sit
es_rarefaction/main/data/example_vcf.vcf"),
                   header = TRUE)

# download the vcf (if copy/pasting, check line breaks)
download.file("https://raw.githubusercontent.com/ChristieLab/seg_sites_rarefa
ction/main/data/example_vcf.vcf",
              destfile = "example_vcf.vcf")
```

32 This VCF file contains a subset of data from Hemstrom et al. (W. B. Hemstrom et al., 2022)

33 which includes gennotypes for 1,000 SNP loci from five populations. The VCF file can be

34 loaded in alongside the metadata into a single object using "read_vcf":

```r
monarchs <- read_vcf("example_vcf.vcf", sample.meta = meta)
```

36 Population size information can be accessed using "summarize_facets" by referring to the

37 column name in the metadata read in earlier ("pop"). "Facets" in snpR refer to any metadata

38 column in the data (including both sample metadata, like we read in above, and locus metadata if

39 supplied). We can view the number of individuals per population using "summarize_facets":

```r
summarize_facets(monarchs, "pop")
```

41 Running this will show the number of samples per population.

42        The expected number of segregating sites per population can be calculated using the

43    function "`calc_seg_sites`", supplying the object we imported above and naming the facet

44    which contains our population information. The rarefaction level ($\gamma$) per locus will be

45    automatically calculated according to the argument "*g*". If "*g*" is zero (the default), $\gamma = n_{min}$,

46    where $n_{min}$ is the smallest sample size across all populations for each locus after accounting for

47    missing data. On the other hand, if g < 0, $\gamma$ will be set to $n_{min} - g$ and if g > 0, $\gamma$ will be set to g.

48    Either way, the result can be fetched with "`get.snpR.stats`", referring to the object, facet, and
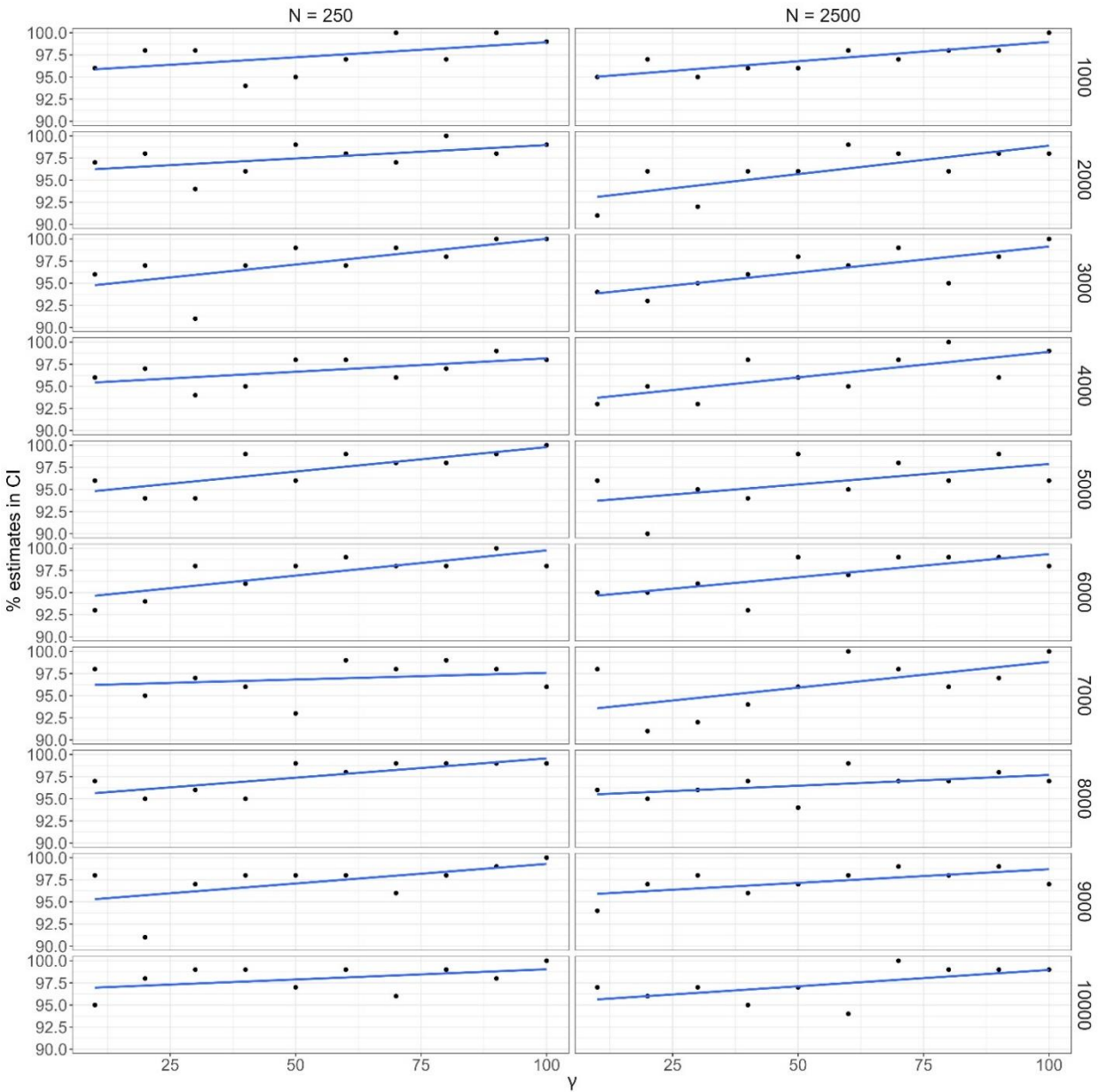
49    statistic we are fetching:

```
50   # g = 0, gamma = Nmin
51   monarchs <- calc_seg_sites(monarchs, "pop", g = 0)
52   get.snpR.stats(monarchs, "pop", "seg_sites")$weighted.means
53
54   # g = -1, gamma = Nmin - 1
55   monarchs <- calc_seg_sites(monarchs, "pop", g = -1)
56   get.snpR.stats(monarchs, "pop", "seg_sites")$weighted.means
57
58   # g = 10, gamma = 10
59   monarchs <- calc_seg_sites(monarchs, "pop", g = 10)
60   get.snpR.stats(monarchs, "pop", "seg_sites")$weighted.means
```

61    Note that the addition of "`$weighted.means`" to the end of each "`get.snpR.stats`" means that

62    we are fetching the mean values specifically, not the per-locus data. We can fetch that instead by

63    using "`$single`", referring to the statistics for each single locus. The mean results will contain

64    the columns "`seg_sites`" and "seg_sites_var" containing $E(N_S)$ and its variance $\sigma^2_{E(N_S)}$,

65    respectively. The per-locus results will contain the columns "g_prob_seg", "prob_seg", and

66    "prob_seg_var" which note $\gamma$, the probability the site segregates at $\gamma$ in a specific population

67    ($\hat{P}(S_{ij})$), and the variance of $\hat{P}(S_{ij})$, $\sigma^2_{\hat{P}(S_{ij})}$, respectively.
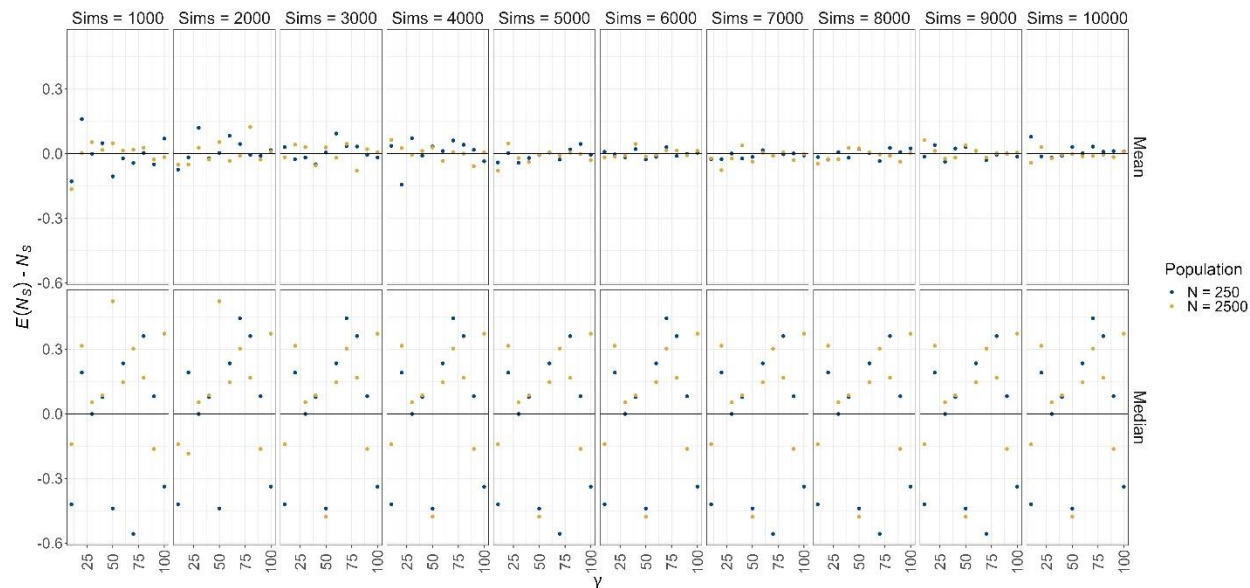
68 **Supplementary Figures:**



69

70 **Figure S1:** The percentage of expected segregation probabilities $\hat{P}(S_{iq})$ for which the

71 probability a loci was segregating was inside the 95% confidence interval derived from

72 simulations for different γ values and numbers of simulations for both a population of size 250
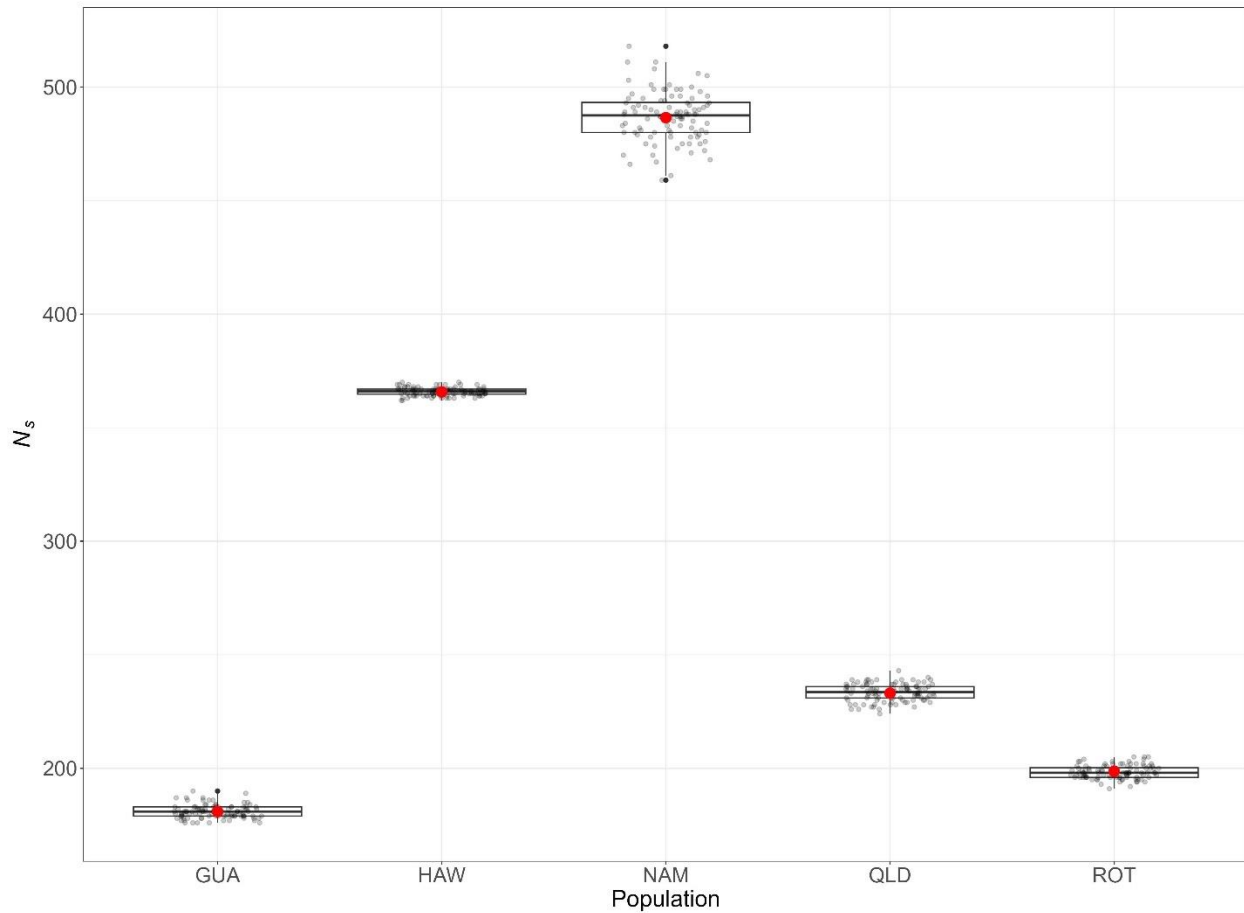
73 and 2,500.

74



75

**Figure S2:** The difference between the expected number of segregating sites ($E(N_S)$) and the

mean and median observed number from simulations across different numbers of

simulations and γ values for populations of both size 250 and 2,500.

79

**Figure S3:** The expected number of segregating sites ($E(N_S)$, red) compared to the distributions

of the number of segregating sites observed after rarefaction based on 100 simulations (black).

$E(N_S)$ values were obtained using the "calc_seg_sites" function in snpR using "g = 0", which

sets γ equal to the minimum sample size across all populations for each locus.