

Diversity and overlap analysis in TCR populations

Michał Seweryn and Grzegorz A Rempala
mseweryn@cph.osu.edu, grempala@cph.osu.edu



Problem

Maintaining a proper diversity of T-cell receptor populations (TCR's) is crucial for the immune system's ability to recognize a vast variety of foreign antigens and to avoid autoaggression. Due to large diversity of TCR's and the sampling error (of high-throughput sequencing) the standard diversity and overlap measures of the contingency table analysis are insufficient. Applying information theory here we have developed some new ones specifically for TCR data analysis.

Basic Concepts

Each population of TCR's corresponds to a vector of counts $\mathbf{c}_i = (c_{1,i}, \dots, c_{m,i}, i), i = 1, \dots, s$. Let $X = (X_1, \dots, X_m), \sum_{k=1}^m X_k = n$ be a sample of size n . We define 'sample coverage' $C := \sum_{l=1}^m p_l I_l$, where $I_l = 1$ if $X_l > 0$ and $I_l = 0$ otherwise and it's Good and Turing estimator $\hat{C} = \frac{s_n(1)}{n}$, where $s_n(1)$ is the number of singletons.

Let \mathcal{F}_c be the fingerprint or diversity of the population \mathbf{c} – a vector given by $\mathcal{F}_c = (s(1), \dots, s(\max_i c_i))$, $s(k) = \text{card}\{i : c_i = k\}$. A nonnegative, real function of the fingerprint is called a measure (index) of diversity.

Set $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n$ to be populations and let $\text{supp}(\mathbf{c}_i)$ denote the support of \mathbf{c}_i . The overlap between vectors $\mathbf{c}_1, \dots, \mathbf{c}_n$ is then $\mathcal{O}_n = \bigcap_{k=1}^n \text{supp}(\mathbf{c}_k)$. Any nonnegative, real function G , such that $G(\mathcal{O}_n \mathbf{c}_1, \dots, \mathbf{c}_n)$, is an overlap measure (index).

Let \mathbf{c} - population and D - monotone diversity measure, $\mathcal{F}_{I(m)}$ the fingerprint of a uniform population with m different receptors. We define ENS (effective number of species) as the smallest solution of the equation $D(\mathcal{F}_{I(y)}) = D(\mathcal{F}_c)$.

Previous work – diversity and overlap indices

Renyi's (and Shannon's) entropy $H_\alpha(\mathcal{F}_c) = \frac{1}{1-\alpha} \log \left(\sum_k s(k) \left(\frac{k}{n} \right)^\alpha \right)$, $\alpha \geq 0$, and $H_1(\mathcal{F}_c)$

Simpson's index $ISI := \exp(H_2(\mathcal{F}_c))$, Chao-Shen's index $H_1(\mathcal{F}_c) = -\sum_k s(k) \frac{k}{n(1-(k/n)^n)}$

For two populations \mathbf{c}, \mathbf{c}' define Jaccard's index $J(\mathbf{c}, \mathbf{c}') = \frac{\sum_i \min(c_i, c'_i)}{\sum_i (c_i + c'_i) - \sum_i \min(c_i, c'_i)}$

For $(\mathbf{p}_1, \mathbf{p}_2)$ – pair of normalized populations: Morisita-Horn's index and Renyi's divergence

$$MH(\mathbf{p}_1, \mathbf{p}_2) = \frac{2 \sum_i p_{i,1} p_{i,2}}{\sum_i p_{i,1}^2 + \sum_i p_{i,2}^2}, \quad F_\alpha(\mathbf{p}_1, \mathbf{p}_2) = \frac{1}{\alpha-1} \log \left(\sum_i \frac{p_{i,1}^\alpha}{p_{i,2}^\alpha} \right), \alpha \geq 0$$

Coverage corrected diversity indices

Now we aim to estimate the diversity of a population of T-cells given sample \mathcal{X} of size n . We define a family of 'sample based' diversity measures which in a natural way overemphasize rare species in the case of undersampling error, we also use a Horvitz-Thompson type correction for bias. Let \hat{C} be the Good-Turing estimator of the sample coverage, we define

$$\hat{H}_{\alpha \hat{C}}(\mathcal{X}) = \frac{\log \left(\sum_k s_n(k) \left(\frac{k}{n} \right)^{\alpha \hat{C}} \right)}{1 - \alpha \hat{C}}, \quad \hat{H}_{\alpha \hat{C}}^{(n)}(\mathcal{X}) = \frac{\log \left(\sum_k \frac{s_n(k) k^{\alpha \hat{C}}}{n^{\alpha \hat{C}} (1 - (k/n)^n)} \right)}{1 - \alpha \hat{C}}$$

Set $\hat{\mathbf{p}}$ to be the MLE of the normalized population vector \mathbf{p} and $\hat{\mathbf{p}} = \hat{C} \hat{\mathbf{p}}$. Let $0 < \alpha < \infty$ and assume that $H_\alpha(\mathbf{p}) < \infty$. If $\alpha < 1$ or if $\alpha > 1$ and $\sum_k p_k \log^r(1/p_k) < \infty$ for some $r > 0$ then

$$H_\alpha^{(n)}(\hat{\mathbf{p}}) \xrightarrow{a.s.} H_\alpha(\mathbf{p}) \quad \text{and} \quad H_{\hat{C}\alpha}^{(n)}(\hat{\mathbf{p}}) \xrightarrow{a.s.} H_\alpha(\mathbf{p}). \quad \text{If } \alpha = 1, \text{ then } H_1^{(n)}(\hat{\mathbf{p}}) \xrightarrow{a.s.} H_1(\mathbf{p}),$$

and on the set $\{\hat{C} < 1 \text{ i.o.}\}$, $H_{\hat{C}}^{(n)}(\hat{\mathbf{p}}) \xrightarrow{\text{in prob.}} H_1(\mathbf{p})$, where $R_1^{(n)}(\hat{\mathbf{p}}) := \sum \frac{\hat{p}_i}{1 - (1 - \hat{p}_i)^n}$

New overlap indices

We consider a slightly more general form of the Morisita-Horn index, which allows it to put more weight on rare (resp. abundant) receptors. For $(\mathbf{p}_1, \mathbf{p}_2)$ a pair of normalized populations and $\alpha, \beta \in (0, \infty)$ the power-geometric (or PG) index of overlap is given by

$$PG_{\alpha, \beta}(\mathbf{p}_1, \mathbf{p}_2) = \frac{\sum p_{i1}^\alpha p_{i2}^\beta}{\sum p_{i1}^{2\alpha} + \sum p_{i2}^{2\beta}}$$

In analogy with the adjustment of diversity indices, and in the notation as above, we may consider $PG_{\hat{C}_1 \alpha, \hat{C}_2 \beta}^{(n)}(\hat{\mathbf{p}}_1, \hat{\mathbf{p}}_2)$ as the sample-coverage and Horvitz-Thompson adjusted PG index.

Assume that $\sum p_{i1}^\alpha < \infty$ and $\sum p_{i2}^\beta < \infty$, as well as $\sum p_{i1} \log^{r_1} \frac{1}{p_{i1}} < \infty$ for some $r_1 > 0$, if $\alpha > 1$ and $\sum p_{i2} \log^{r_2} \frac{1}{p_{i2}} < \infty$ for some $r_2 > 0$, if $\beta > 1$. Then

$$PG_{\hat{C}_1 \alpha, \hat{C}_2 \beta}^{(n)}(\hat{\mathbf{p}}_1, \hat{\mathbf{p}}_2) \xrightarrow{a.s.} PG_{\alpha, \beta}(\mathbf{p}_1, \mathbf{p}_2).$$

Moreover, we consider a different approach based on $(m \times n)$ contingency table $\mathbf{C} = [c_{ij}]$ with columns representing n different population and rows representing m receptors. Let $\mathbf{P} = [p_{ij}] := [\frac{c_{ij}}{\sum_{kl} c_{kl}}]$ be the normalized matrix with columns $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n$, $\mathbf{p}_{i\circ} = \sum_j p_{ij}$, $\mathbf{p}_{\circ j} = \sum_i p_{ij}$ and $\mathbf{P}_\circ = (p_{\circ 1}, \dots, p_{\circ n})$, $\mathbf{P}^\circ = (p_{1\circ}, \dots, p_{m\circ})$, as well as $\mathbf{Q} = \mathbf{P}_\circ \otimes \mathbf{P}^\circ := [p_{i\circ} p_{\circ j}]$. Define

$$I_\alpha(\mathbf{C}) = 1 - F_\alpha(\mathbf{P}, \mathbf{Q}) / H_{2-\alpha}(\mathbf{P}_\circ) \quad \text{and} \quad Q_\alpha(\mathbf{C}) = 1 - I_\alpha(\mathbf{C}).$$

Note that for $\alpha \in (0, 2)$ we have $0 \leq Q_\alpha(\mathbf{C}) \leq 1$, $Q_\alpha(\mathbf{C}) = 0$ iff $\mathbf{p}_1 = \mathbf{p}_2 = \dots = \mathbf{p}_n$ and if the vectors $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n$ form an orthogonal system, then $Q_\alpha(\mathbf{C}) = 1$.

Let $\hat{\mathbf{P}}$ be the empirical MLE of \mathbf{P} , then we also have that $I_\alpha(\hat{\mathbf{P}}) \xrightarrow{a.s.} I_\alpha(\mathbf{P})$.

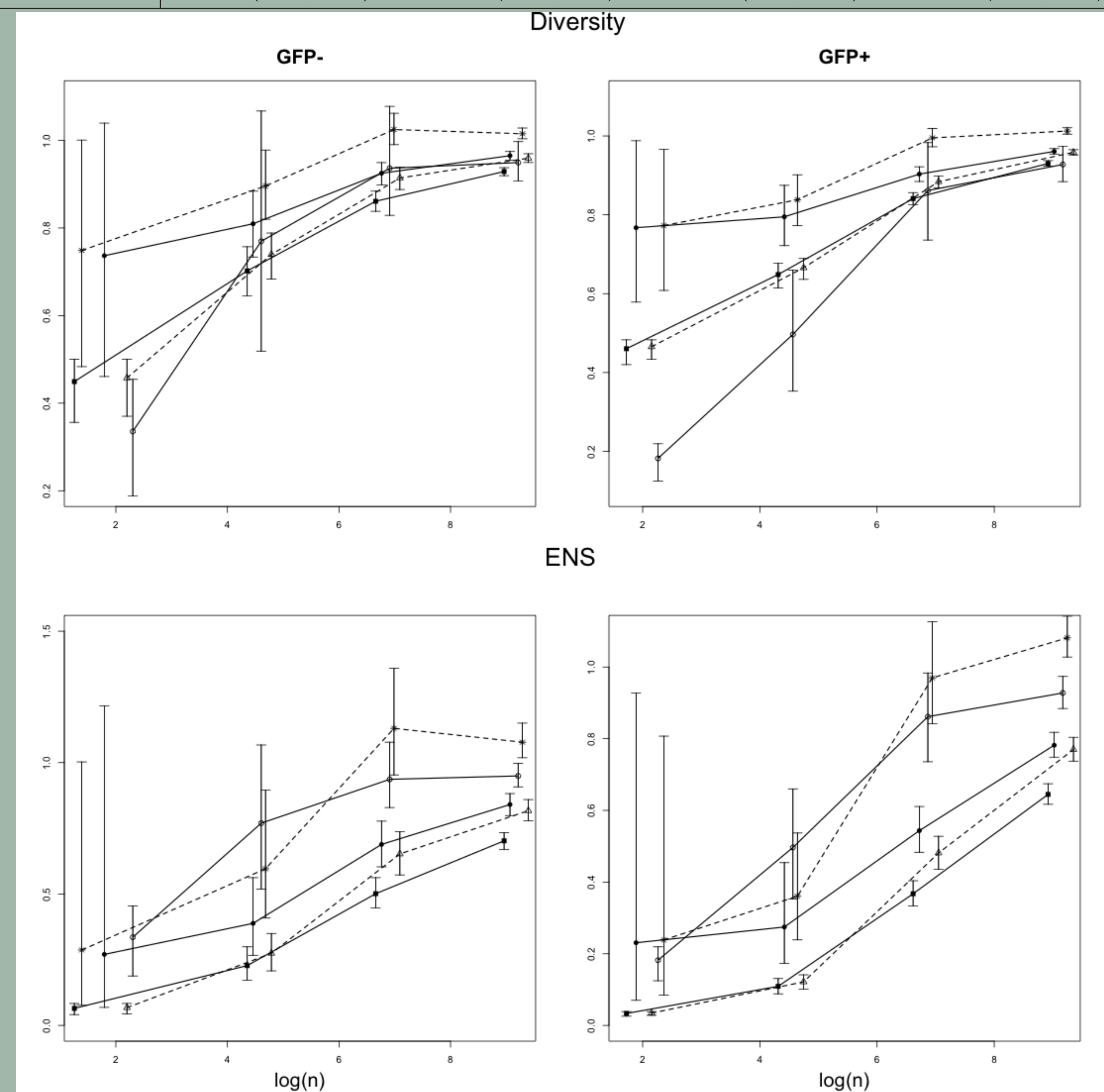
Acknowledgements

This research was partially supported by US NIH grant R01CA-152158 (GAR, MS) and US NSF grant DMS-1106485 (GAR). The supports are gratefully acknowledged.

Performance of new diversity indices

We analyze two TCR datasets obtained from high-throughput sequencing experiments conducted in the molecular immunology lab of Prof Leszek Ignatowicz. One dataset consists of the so called "regulatory" T-cells (GFP^+) the second one of the so-called "naive" T-cells. Diversity and ENS (based on 500 repetitions) is reported relatively to the values in the complete set.

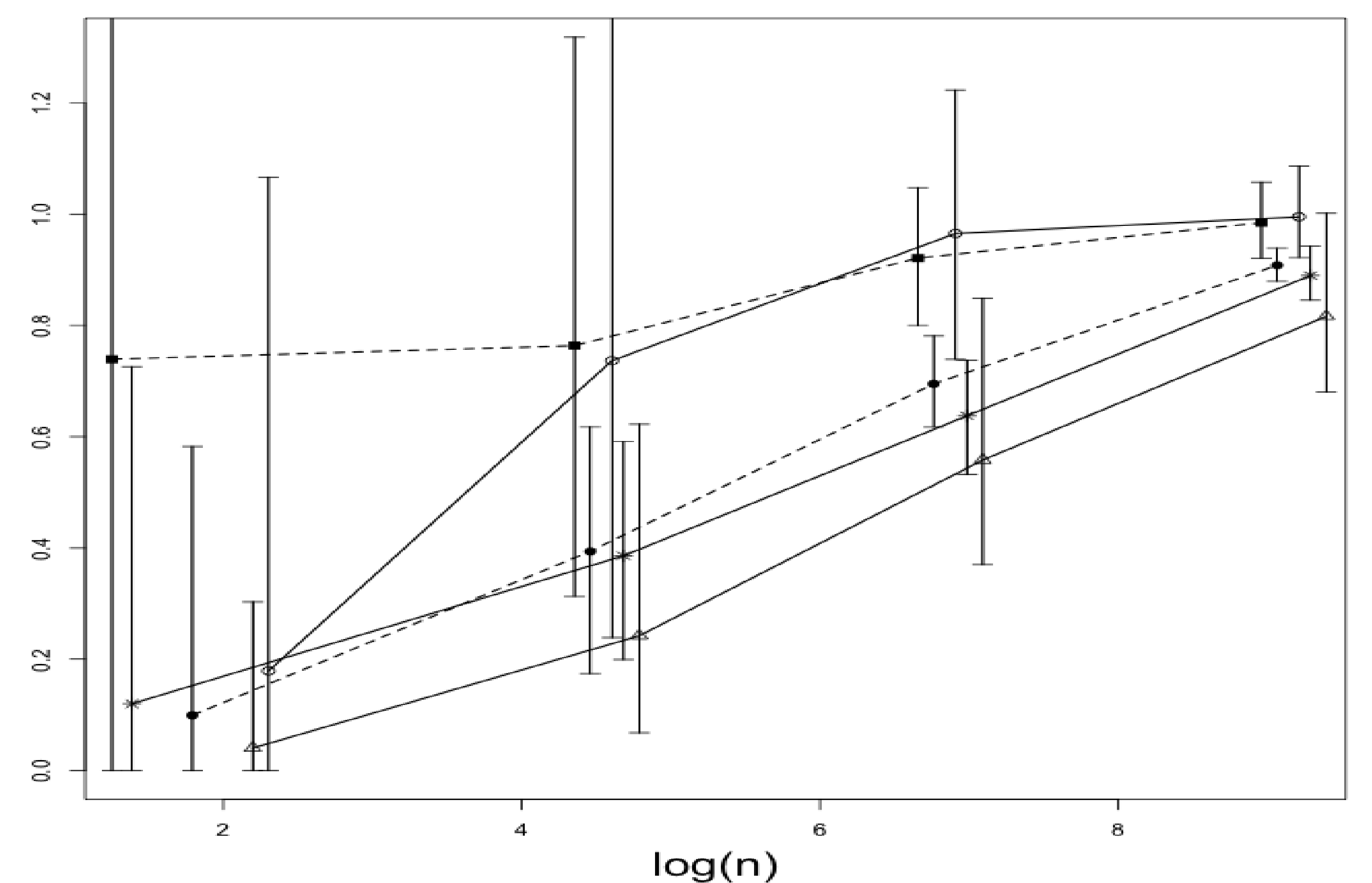
Stat/ENS	$n = 10^2$ $\hat{C} = 0.30$	$n = 10^3$ $\hat{C} = 0.62$	$n = 10^4$ $\hat{C} = 0.83$	$n = 10^5$ $\hat{C} = 0.94$
ISI	0.34 (0.19,0.45) 0.34 (0.19,0.45)	0.77 (0.52,1.07) 0.77 (0.52,1.07)	0.94 (0.83,1.08) 0.94 (0.83,1.08)	0.95 (0.90,0.99) 0.95 (0.90,0.99)
$H_{\hat{C}}$	0.46 (0.37,0.50) 0.07 (0.04,0.08)	0.74 (0.68,0.78) 0.27 (0.21,0.35)	0.92 (0.89,0.94) 0.65 (0.57,0.74)	0.96 (0.95,0.97) 0.83 (0.78,0.86)
$H_{\hat{C}}^{(n)}$	0.75 (0.49,1.00) 0.29 (0.077,1.00)	0.90 (0.82,0.98) 0.60 (0.41,0.90)	1.02 (1.00,1.06) 1.13 (0.95,1.35)	1.01 (1.00,1.02) 1.07 (1.01,1.15)
$H_1^{(n)}$	0.73 (0.46,1.04) 0.27 (0.06,1.22)	0.80 (0.73,0.89) 0.39 (0.26,0.56)	0.92 (0.90,0.95) 0.69 (0.60,0.78)	0.96 (0.95,0.97) 0.84 (0.79,0.88)
Plug-in H_1	0.45 (0.36,0.50) 0.06 (0.04,0.08)	0.70 (0.65,0.76) 0.23 (0.17,0.30)	0.86 (0.84,0.88) 0.50 (0.44,0.56)	0.93 (0.92,0.94) 0.70 (0.67,0.73)



Solid lines, from the top: (i) ISI , (ii) Plug-in, (iii) $H_1^{(n)}$; dashed lines: (i) $H_{\hat{C}}$, and (ii) $H_{\hat{C}}^{(n)}$.

Performance of new overlap indices

Stat	$n = 10^2$ $\hat{C}_1 = 0.25$ $\hat{C}_2 = 0.16$	$n = 10^3$ $\hat{C}_1 = 0.61$ $\hat{C}_2 = 0.40$	$n = 10^5$ $\hat{C}_1 = 0.83$ $\hat{C}_2 = 0.70$	$n = 10^6$ $\hat{C}_1 = 0.94$ $\hat{C}_2 = 0.91$
PG	0.74 (0.00,4.2)	0.76 (0.31,1.31)	0.92 (0.80,1.04)	0.99 (0.92,1.05)
I_1 -ind	0.10 (0.00,0.59)	0.40 (0.18,0.62)	0.69 (0.62,0.78)	0.91 (0.88,0.95)
L	0.12 (0.00,0.73)	0.38 (0.20,0.59)	0.64 (0.53,0.74)	0.88 (0.84,0.94)
CJ	0.04 (0.00,0.30)	0.24 (0.06,0.62)	0.56 (0.37,0.85)	0.81 (0.68,1.01)
MH	0.17 (0.00,1.07)	0.74 (0.23,1.43)	0.96 (0.73,1.22)	0.99 (0.92,1.09)



Solid lines: (i) MH (open circles), (ii) L (stars) and (iii) CJ (triangles); dashed lines (i) $PG_{\hat{C}_1, \hat{C}_2}^{(n)}$ (squares), and I -index (filled circles).