1 **Supplementary information**

2

3 <u>Materials and methods: 454-pyrosequencing and noise removal</u>

4 Tag-pyrosequencing was done with Roche 454 Titanium platform following

5 manufacturer protocols (454 Life Science). Amplification of the hypervariable regions

6 V1-V3 was done using Primers 28F (5'-GAGTTTGATCNTGGCTCAG) and 519R (5'-

7 GTNTTACNGCGGCKGCTG). Approximately 400 bp long tags were obtained. PCR

8 and subsequent sequencing are described in Dowd *et al.* (2008).

9     The raw tag-sequences were processed using QIIME (Caporaso *et al.*, 2010). Briefly,

10 to reduce sequencing errors and their effects, the multiplexed reads were first trimmed,

11 quality-filtered and assigned to the samples, surface or bottom. The filtering criteria

12 included a perfect match to the sequence barcode and primer, at least 400 bp in length, an

13 average quality score (phred) of 28 within sliding windows of 50bp. Additionally,

14 denoiser was used to reduce the amount of erroneous sequences (Reeder & Knight,

15 2010).  The sequences were then clustered into Operational Taxonomic Units (OTUs)

16 based on the relatedness of the sequences (97% similarity) with UCLUST, version

17 1.1.579 (Edgar, 2010). A representative sequence from each OTU was selected as the

18 first cluster seed chosen by UCLUST.  ChimeraSlayer implemented in Mothur (Schloss

19 *et al.*, 2011) use to check for chimeras. Then, taxonomy assignment was made with

20 QIIME by searching the representative sequences of each OTU against the SILVA

21 16S/18S rDNA non-redundant reference dataset (SSU Ref 108 NR) (Quast *et al.*, 2013)

22 using the Basic Local Alignment Search Tool (BLAST) and an e-value of 0.03. Chimera,

23 chloroplast, eukarya and archaea sequences were removed from the output fasta file that

24  was used for building a table with the OTU abundance of each sample and the taxonomic

25  assignments for each OTU.

26

27  Materials and methods: Isolation of bacterial cultures

28  Isolates were obtained on board by plating 100 µl of undiluted and 10x diluted sea-water

29  from the surface sample, in triplicates, onto modified Zobell agar plates (i.e. 5 g peptone,

30  1 g yeast extract and 15 g agar in 1 l of 0.2 µm filtered 75% sea water). Agar plates were

31  incubated at *in situ* temperature (~20 °C), in the dark, for 14 days. 326 bacterial colonies

32  were selected and the cultures were subsequently purified by re-isolation three times in a

33  month. Next, isolates were grown at 20 ºC on the same liquid medium and stored at -

34  80 ºC with 25% (v/v) glycerol.

35

36  Materials and methods: Bacterial isolates PCRs

37   PCR, using Taq polymerase (Boehringer-Mannheim), of the Internal Transcribed Spacer

38  (ITS) was done using primers ITS-F (5'-GTCGTAACAAGGTAGCCGTA) and ITS-R

39  (5'-GCCAAGGCATCCACC) and the following thermal conditions: 94ºC for 2 min, then

40  32 cycles of 94 ºC for 15 sec, 55 ºC for 30 sec, 72 ºC for 3 min, followed by one cycle of

41  72 ºC for 4 min and 4 ºC on hold.

42   PCR of the 16S rRNA gene of the 148 chosen by their different ITS pattern were then

43  amplified using bacterial 16S rRNA gene primers 27F

44  (5'-AGAGTTTGATCMTGGCTCAG) and 1492R (5'-GTTTACCTTGTTACGACTT).

45  The thermal conditions were as follows: 94 ºC for 5 min, then 30 cycles of 94 ºC for 1

46    min, 55 ºC for 1 min, 72 ºC for 2 min, followed by one cycle of 72 ºC for 10 min and 4 ºC

47    on hold.

48

49    <u>Materials and methods: Simulating the Required Sequencing Effort (RSE)</u>

50    For each of 80 ensemble members, we simulated a random sequence of 10N individual

51    species labels, where N is the present sequencing effort, by: 1) sampling a set of

52    parameter values from the posterior distribution, 2) sampling relative abundances

53    (proportions in the water sample) from the taxon abundance distribution given the

54    parameter values, 3) sampling species counts (from hypothetical sequencing) using the

55    multinomial distribution given the proportions and the total number of individuals 10N,

56    and 4) converting the species counts into a randomly-ordered sequence of individual

57    labels. The simulated RSE was then identified as the individual (tag) index for which the

58    number of species observed earlier in the sequence first exceeded 90% of the simulated

59    total richness (S).

60

61    <u>Discussion: Simulation tests on the number of isolates retrieved in sequences</u>

62    To simulate the number of isolates retrieved in sequences, we simulated 3000 sets of

63    species counts using the method described above for RSE calculations, but with the total

64    number of tags fixed at the present sequencing effort N.  For each set of simulated

65    sequencing counts, 38 species were selected at random without replacement from the list

66    of all S counts (including zeros), and the number of these with non-zero counts was

67    recorded to give the simulated number retrieved by sequencing $r_s$. The simulation p-

68      value for the actual number of species retrieved by sequencing $r$ was then taken as $(1 +$

69      $\#(r_s \leq r))/3001$ following Davison & Hinkley (1997).

70

71      <u>Discussion: Simulation/bootstrap tests on the counts of isolates retrieved in sequences</u>

72      To simulate the counts of isolates retrieved in sequences, we again simulated 3000 sets of

73      species counts as described above, and this time randomly selected without replacement

74      either 9 or 14 species from the list of non-zero counts for each simulation.  The mean,

75      median, and maximum counts from this subset were recorded for each simulation, and p-

76      values were calculated as described above assuming lower-than-random count statistics

77      as an alternative hypothesis.  For the bottom count statistics (14 species), we also

78      performed the test assuming higher-than-random count statistics $t$ on the alternative,

79      hence calculating p-values as $(1 + \#(t_s \geq t))/3001$.

80          These tests were also repeated using a bootstrap method, thus avoiding the need to

81      assume a parametric distribution.   To do this, a vector of 9 or 14 species counts was

82      randomly resampled *with* replacement from the observed species count vector.  This was

83      repeated over 9999 bootstraps and bootstrap p-values were calculated as $(1 + \#(t_s \leq$

84      $t))/10000$ or $(1 + \#(t_s \geq t))/10000$, again following Davison & Hinkley (1997).

85

86

87

88

89

90

91    References:

92
93    Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman F, Costello E, *et al.*
94        (2010). QIIME allows analysis of high- throughput community sequencing data.
95        *Nature* **7**:335–336.
96    Davison A, Hinkley D. (1997). Bootstrap methods and their applications. Cambridge
97        University Press: New York.
98    Dowd SE, Callaway TR, Wolcott RD, Sun Y, McKeehan T, Hagevoort RG, *et al.* (2008).
99        Evaluation of the bacterial diversity in the feces of cattle using 16S rDNA bacterial
100       tag-encoded FLX amplicon pyrosequencing (bTEFAP). *BMC Microbiol* **8**:125.
101   Edgar RC. (2010). Search and clustering orders of magnitude faster than BLAST.
102       *Bioinformatics* **26**:2460–2461.
103   Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, *et al.* (2013). The SILVA
104       ribosomal RNA gene database project: improved data processing and web-based tools.
105       *Nucleic Acids Res* **41**:590–596.
106   Reeder J, Knight R. (2010). Rapidly denoising pyrosequencing amplicon reads by
107       exploiting rank-abundance distributions. *Nat Methods* **7**:668-669.
108   Schloss PD, Gevers D, Westcott SL. (2011). Reducing the effects of PCR amplification
109       and sequencing artifacts on 16S rRNA-based studies. Gilbert, JA (ed). *PLoS One*
110       **6**:e27310.
111

112

113

114

115

116    **Figure legends**

117    Figure S1.  Goodness-of-fit of the best-approximating Sichel distribution to surface (A)

118    and bottom (B) data sets.  Observed and predicted count frequencies (numbers of OTUs

119    with a given sample abundance) are plotted against tag counts (sample abundances) on a

120    log-log scale.  Goodness-of-fit is illustrated by the closeness of the predicted frequencies

121    (posterior means, solid lines) to the observed frequencies (dots) as well as by the

122    narrowness of the 95% prediction intervals (dashed lines) while still containing most of

123    the data.  The comparison is restricted to rare counts in the range 1–100 because these are

124    likely the most important for estimating total richness and required sequencing effort, and

125    because the computation of stable frequency prediction intervals for higher counts would

126    require too many simulations (the intervals shown used 3 000).  The distributions were

127    however fitted to the full range of observed count frequencies ($f_{1-178569}$ and $f_{1-45414}$ for

128    surface and bottom samples respectively).

129

130

131    **Table captions**

132    Table S1.  Four different compound Poisson distributions were fitted to the surface and

133    bottom data: the Poisson log-normal, the Poisson inverse Gaussian, the Poisson log-

134    student, and the Poisson generalized inverse Gaussian (Sichel) distribution.  As a

135    robustness check we reran the Sichel fit for the surface sample excluding the counts of

136    the most abundant species which, for this sample, was more than 3 times as abundant as

137    the second most abundant species (see Surface*).  The relative goodness-of-fit is assessed

138    using Akaike's Information Criterion (AICc = -2×max(log likelihood) + $2p$ + $2p(p+1)/(n-$

139    $p-1)$, where $p$ is the number of fitted parameters and $n$ is the number of data; Hurvich and

140    Tsai, 1989; Burnham and Anderson, 2002) and the deviance information criterion (DIC =

141    -2×posterior mean(log likelihood) + $p$; Spiegelhalter *et al*., 2002; Quince *et al*., 2008).

142    For the robustness check the selection criteria are placed in square parentheses since these

143    cannot be compared to other rows.  We also show the total species richness estimates

144    from maximum likelihood ($\hat{S}_{ML}$) as well as the posterior median ($\hat{S}_{50\%}$) and the 95%

145    credible bounds ($\hat{S}_{2.5\%}$ and $\hat{S}_{97.5\%}$) from the Bayesian MCMC method (Quince *et al*.

146    2008).

147    Reference: Hurvich CM and Tsai C-L (1989). Regression and time series model selection in

148    small samples. *Biometrika* **76**: 297-307.

149

150    Table S2. (A) Semiparametric functional fits to surface sample collector's curve data and

151    corresponding estimates of total species richness.  A set of 12 convex, saturating

152    functions were fitted to the rarefied species accumulation curve, sampled at intervals of 1

153    000 tags (hence 502 data points), using the nonlinear least squares function "nls" in R to

154    estimate the parameters a, b etc.  The absolute quality of the fits was measured using the

155    generalized R2 values (defined for nonlinear fit as 1 - RSS/SSM, where RSS is the

156    residual sum of squares and SSM is the sum of squares of the sample mean).  The best-

157    approximating model was selected as that which minimized Akaike's Information

158    Criterion (AICc, in this case the Power Michaelis Menten (2) function was selected).  The

159    selected model was then used to estimate the total sample richness S as the asymptotic

160    value of the function at x = Inf (final column shows the estimates for all candidate

161    functions).  Required sequencing effort (not shown) was predicted by inverting the

162    selected function for x such that the value of the function was 0.9 times the estimated

163    sample richness.  Note that for certain 3 and 4 -parameter functions the R2 values are

164    extremely high and differ only in the fourth or fifth decimal places (R2>0.999) yet the

165    estimated richnesses can differ substantially (cf. Power Michaelis Menten (2) vs. Weibull

166    Cumulative).  For such functions, the AICc values also tend to differ by relatively large

167    amounts, such that a model averaging strategy based on AIC weights would differ little

168    from simply choosing the lowest-AICc model (Burnham and Anderson, 2002), and any

169    assessment of model selection uncertainty based on AIC-weights is unlikely to predict the

170    level of selection uncertainty observed in simulations (see Table S3).  This latter is likely

171    the result of the neglected error correlation in the functional fits.

172    Table S2. (B) Semiparametric functional fits to bottom sample collector's curve data and

173    corresponding estimates of total species richness.  A set of 12 convex, saturating

174    functions were fitted to the rarefied species accumulation curve, sampled at intervals of 1

175    000 tags (hence 576 data points), using the nonlinear least squares function "nls" in R to

176    estimate the parameters a, b etc.  The absolute quality of the fits was measured using the

177    generalized R2 values (defined for nonlinear fit as 1 - RSS/SSM, where RSS is the

178    residual sum of squares and SSM is the sum of squares of the sample mean).  The best-

179    approximating model was selected as that which minimized Akaike's Information

180    Criterion (AICc, in this case the Power Michaelis Menten (2) + offset function was

181    selected).  The selected model was then used to estimate the total sample richness S as the

182    asymptotic value of the function at x = Inf (final column shows the estimates for all

183    candidate functions).  Required sequencing effort (not shown) was predicted by inverting

184    the selected function for x such that the value of the function was 0.9 times the estimated

185    sample richness.  Note that for certain 3 and 4 -parameter functions the R2 values are

186    extremely high and differ only in the fourth or fifth decimal places (R2>0.999) yet the

187 estimated richnesses can differ substantially (cf. Power Michaelis Menten (2) vs. Weibull

188 Cumulative). For such functions, the AICc values also tend to differ by relatively large

189 amounts, such that a model averaging strategy based on AIC weights would differ little

190 from simply choosing the lowest-AICc model (Burnham and Anderson, 2002), and any

191 assessment of model selection uncertainty based on AIC-weights is unlikely to predict the

192 level of selection uncertainty observed in simulations (see Table S3). This latter is likely

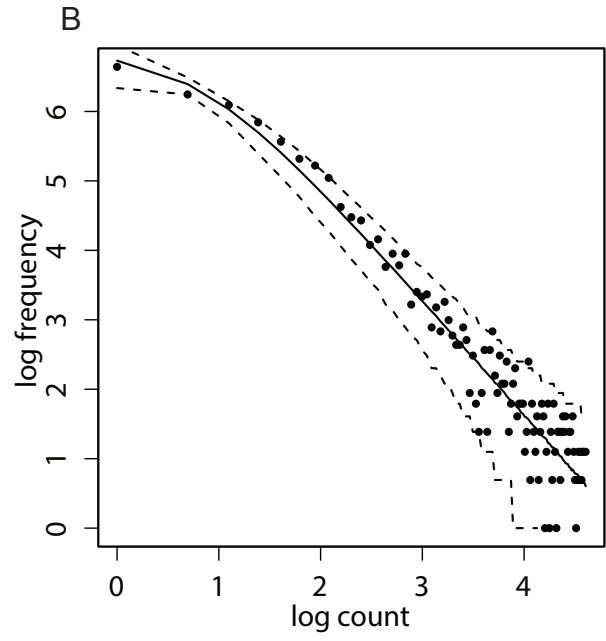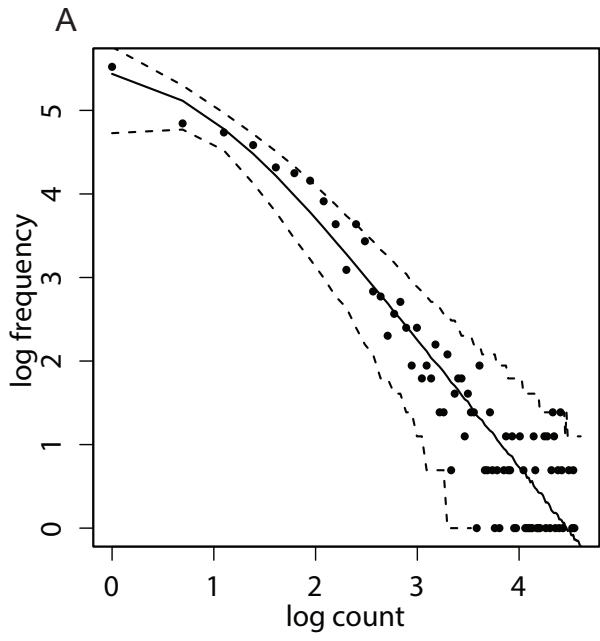193 the result of the neglected error correlation in the functional fits.

194

195 Table S3. Simulation-based tests of estimator performance, considering estimates of both

196 the total species richness (S) and the required sequencing effort (RSE) i.e. number of tags

197 required to observe a given fraction of the total richness in a new sample (e.g. 0.7S means

198 70% of the total richness). For each of four parametric distributions (Poisson log-normal,

199 Poisson log-student, Poisson inverse-Gaussian, and Sichel) an ensemble of 80 sets of

200 community abundances were randomly sampled from the parametric distribution; species

201 data were then simulated by sampling from multinomial distributions with probabilities

202 defined by the community abundances for each ensemble member. Distribution

203 parameter values, including the total species richness, were also varied between ensemble

204 members by sampling from the posterior distributions fitted to the observed data.

205 Estimator performance is summarized by the %BIAS (ensemble average of estimate

206 minus true value) and %RMSE (root-mean-square error), normalizing by the ensemble

207 mean of the true value in both cases. Non-parametric species richness estimators

208 included the Chao1 lower bound estimate (Chao, 1984), the coverage-based estimator for

209 highly heterogeneous communities (ACE-1; Chao & Lee, 1992; Chao et al., 2000) and

210    the bias-corrected Chao estimate iChao (Chiu *et al*., 2014).  The ACE-1 estimator was

211    tested using two values of the cut-off count k to define "rare" species: the default value k

212    = 10 and a larger value k = 100 as recommended by Chao & Shen (2012) for microbial

213    communities (note, the estimated CV of the "rare" species was < 0.8 for k = 10 but > 0.8

214    for k = 100, where 0.8 is a threshold above which Chao & Shen (2012) recommend ACE-

215    1 in preference to ACE).  RSE was estimated for each nonparametric estimator by

216    inverted the expression in Table 1 of Chao *et al*. (2014) and substituting the

217    corresponding estimates of the zero-count frequency f0 = (S - $S_{obs}$) (using ACE-1 this is

218    identical to the method proposed in Chao & Shen (2012) based on Shen *et al*. (2003)

219    except for a negligible bias correction).  Similar results (not shown) were obtained by

220    substituting into equation (12) in Chao *et al*. (2009) (see also Colwell *et al*., 2012,

221    equation 11).  A semi-parametric AICc-selected estimator SP (AICc) was constructed by

222    fitting 12 different functions to the collector's curves (rarefied species richness vs.

223    sampling effort) and choosing the function with the lowest Akaike's Information

224    Criterion (AICc).  Total richness was then estimated as the asymptotic value of the

225    selected function (see Table S2), and RSE was estimated by inverting the selected

226    function for sampling effort given the required fraction of asymptotic richness.

227    Nonparametric estimates were calculated using the R package SPECIES (Wang, 2011)

228    and semiparametric functions were fitted using the nonlinear least squares function "nls"

229    in R (R Core Team, 2013).

230

231

Supplementary Figure 1
Crespo et al.

Table S1.

| Distribution | No. fitted parameters $p$ | Sample | min(-log lik) | AICc | DIC | $\hat{S}_{\text{max. lik.}}$ |
|---|---|---|---|---|---|---|
| Log-normal | 3 | Surface | 869.4 | 1744.8 | 1744.8 | 2449 |
| Log-student | 4 | Surface | 840.6 | 1689.1 | 1689.3 | 1869 |
| Inverse Gaussian | 3 | Surface | 836.9 | 1679.8 | 1679.9 | 1644 |
| Sichel | 4 | Surface | 834.7 | 1677.4 | 1677.3 | 1618 |
| Sichel | 4 | Surface* | [821.3] | [1650.7] | [1651] | 1619 |
| Log-normal | 3 | Bottom | 1276.9 | 2559.8 | 2559.8 | 6843 |
| Log-student | 4 | Bottom | 1198.0 | 2404.1 | 2404.1 | 5850 |
| Inverse Gaussian | 3 | Bottom | 1230.0 | 2466.0 | 2466.0 | 5352 |
| Sichel | 4 | Bottom | 1176.9 | 2361.8 | 2362.1 | 5118 |

| $\hat{S}_{posterior\ mean}$ | $\hat{S}_{50\%}$ | $\hat{S}_{2.5\%}$ | $\hat{S}_{97.5\%}$ |
|---|---|---|---|
| 2501 | 2488 | 2238 | 2819 |
| 1897 | 1891 | 1797 | 2027 |
| 1644 | 1643 | 1594 | 1702 |
| 1615 | 1614 | 1568 | 1669 |
| 1616 | 1615 | 1568 | 1671 |
| 6856 | 6850 | 6544 | 7199 |
| 5867 | 5863 | 5701 | 6055 |
| 5353 | 5352 | 5250 | 5463 |
| 5109 | 5108 | 5027 | 5196 |

Table S2.

A

| Function | Formula (x = #tags-1) | Number of Parameters | $R^2$ | AICc |
|---|---|---|---|---|
| Michaelis Menten | $(ax)/(b+x)+1$ | 2 | 0.98414 | 4976 |
| Negative Exponential | $a(1-\exp(-bx))+1$ | 2 | 0.93681 | 5670 |
| Power Michaelis Menten (1) | $ax^c/(b+x^c)+1$ | 3 | 0.99977 | 2856 |
| Power Michaelis Menten (2) | $ax^c/(b+x)^c+1$ | 3 | 0.99995 | 2086 |
| Power Negative Exponential | $a(1-\exp(-bx))^c+1$ | 3 | 0.99947 | 3274 |
| Weibull Cumulative | $a(1-\exp(-bx)^c)$ | 3 | 0.99992 | 2323 |
| Michaelis Menten + offset | $(ax)/(b+x)+1+c$ | 3 | 0.99680 | 4174 |
| Negative Exponential + offset | $a(1-\exp(-bx))+1+c$ | 3 | 0.98639 | 4901 |
| Power Michaelis Menten (1) + offset | $ax^c/(b+x^c)+1+d$ | 4 | 0.99988 | 2545 |
| Power Michaelis Menten (2) + offset | $ax^c/(b+x)^c+1+d$ | 4 | 0.99995 | 2088 |
| Power Negative Exponential | $a(1-\exp(-bx))^c+1+d$ | 4 | 0.99957 | 3169 |
| Weibull Cumulative + offset | $a(1-\exp(-bx)^c)+d$ | 4 | 0.99992 | 2316 |

B

| Function | Formula | Number of Parameters | $R^2$ | AICc |
|---|---|---|---|---|
| Michaelis Menten | $(ax)/(b+x)+1$ | 2 | 0.98873 | 6899 |
| Negative Exponential | $a(1-\exp(-bx))+1$ | 2 | 0.95179 | 7737 |
| Power Michaelis Menten (1) | $ax^c/(b+x^c)+1$ | 3 | 0.99986 | 4380 |
| Power Michaelis Menten (2) | $ax^c/(b+x)^c+1$ | 3 | 0.99999 | 2897 |
| Power Negative Exponential | $a(1-\exp(-bx))^c+1$ | 3 | 0.99959 | 4996 |
| Weibull Cumulative | $a(1-\exp(-bx)^c)$ | 3 | 0.99999 | 3062 |
| Michaelis Menten + offset | $(ax)/(b+x)+1+c$ | 3 | 0.99758 | 6014 |
| Negative Exponential + offset | $a(1-\exp(-bx))+1+c$ | 3 | 0.98893 | 6891 |
| Power Michaelis Menten (1) + offset | $ax^c/(b+x^c)+1+d$ | 4 | 0.99992 | 4081 |
| Power Michaelis Menten (2) + offset | $ax^c/(b+x)^c+1+d$ | 4 | 0.99999 | 2566 |
| Power Negative Exponential | $a(1-\exp(-bx))^c+1+d$ | 4 | 0.99974 | 4729 |
| Weibull Cumulative + offset | $a(1-\exp(-bx)^c)+d$ | 4 | 0.99999 | 2924 |

$\hat{S}$

1520
1317
1927
1679
1459
1568
1590
1373
1864
1679
1467
1565

$\hat{S}$

4947
4224
6122
5425
4666
4981
5157

4397

5971

5435

4702
4996

Table S3.

| Estimator | Sample | S(lognormal) | | RSE(0.7S, lognormal) | |
|---|---|---|---|---|---|
| | | %BIAS | %RMSE | %BIAS | %RMSE |
| Chao | Surface | -25.8 | 26.8 | -86.2 | 116.1 |
| ACE-1(k=10) | Surface | -24.5 | 25.4 | -85.5 | 114.4 |
| ACE-1(k=100) | Surface | 0.9 | 4.3 | -21.0 | 54.4 |
| iChao | Surface | -23.4 | 24.5 | -84.8 | 114.1 |
| SP(AICc) | Surface | -1.0 | 4.5 | -6.2 | 33.3 |
| Chao | Bottom | -14.9 | 15.1 | -70.1 | 71.9 |
| ACE-1 | Bottom | -14.2 | 14.4 | -70.8 | 72.4 |
| ACE-1(k=100) | Bottom | 7.5 | 8.1 | 54.1 | 55.2 |
| iChao | Bottom | -12.7 | 12.9 | -71.1 | 72.5 |
| SP(AICc) | Bottom | 3.7 | 4.1 | 23.3 | 27.5 |

| S(logstudent) | | RSE(0.8S, logstudent) | | S(inverse Gaussian) | | RSE(0.9S, inv... |
|---|---|---|---|---|---|---|
| %BIAS | %RMSE | %BIAS | %RMSE | %BIAS | %RMSE | %BIAS |
| -22.1 | 22.9 | -76.2 | 90.1 | -5.7 | 6.6 | -36.6 |
| -18.9 | 19.4 | -71.3 | 84.0 | -3.5 | 4.2 | -26.9 |
| 50.2 | 60.0 | 95.2 | 113.4 | 40.6 | 45.9 | 283.0 |
| -19.1 | 19.9 | -72.1 | 85.9 | -3.3 | 4.3 | -25.9 |
| -2.4 | 11.6 | -2.8 | 70.3 | 1.4 | 6.8 | 81.6 |
| -24.3 | 24.8 | -73.9 | 81.3 | -5.9 | 6.0 | -33.5 |
| -18.8 | 19.0 | -66.2 | 72.6 | -3.7 | 3.8 | -23.8 |
| 93.8 | 101.5 | 167.0 | 178.3 | 43.5 | 44.3 | 306.2 |
| -20.6 | 21.0 | -69.0 | 76.4 | -3.2 | 3.5 | -21.6 |
| 3.6 | 10.0 | 31.3 | 60.8 | 1.3 | 7.1 | 51.6 |

| erse Gaussian) | S(Sichel) | | RSE(0.9S, Sichel) | |
|---|---|---|---|---|
| %RMSE | %BIAS | %RMSE | %BIAS | %RMSE |
| 49.6 | -6.3 | 7.8 | -38.3 | 53.9 |
| 36.6 | -3.7 | 4.9 | -27.5 | 40.3 |
| 313.4 | 53.2 | 60.5 | 323.0 | 361.5 |
| 39.2 | -3.6 | 5.5 | -27.5 | 43.8 |
| 448.9 | 4.1 | 12.2 | 237.9 | 921.5 |
| 34.9 | -7.4 | 7.9 | -38.5 | 44.1 |
| 25.0 | -3.6 | 3.9 | -24.7 | 28.5 |
| 309.8 | 76.0 | 81.2 | 396.9 | 416.9 |
| 23.3 | -4.2 | 4.8 | -27.2 | 33.0 |
| 206.8 | 3.1 | 12.5 | 138.5 | 391.6 |