

1 Title:

2 Health and disease imprinted in the time
3 variability of the human microbiome

4 Running title:

5 Microbiota, are you sick?

6 Jose Manuel Martí^{1,2}, Daniel Martínez-Martínez^{1,2,3}, Manuel Peña², César
7 Gracia^{1,2}, Amparo Latorre^{1,3,4,5}, Andrés Moya^{1,3,4,5} & Carlos P. Garay^{1,2,#}

8 ¹Institute for Integrative Systems Biology (I2SysBio), 46980, Spain.

9 ²Instituto de Física Corpuscular, CSIC-UVEG, P.O. 22085, 46071, Valencia, Spain.

10 ³FISABIO, Avda de Catalunya, 21, 46020, Valencia, Spain.

11 ⁴Cavanilles Institute of Biodiversity and Evolutionary Biology, UVEG, 46980, Spain.

12 ⁵CIBER en Epidemiología y Salud Pública (CIBEResp), Madrid, Spain

13 Words count for the Abstract section: 134 of 250 max

14 Words count for the Importance section: 105 of 150 max

15 Words count for the rest of text: 5376 of 5000 max

16

Corresponding author: penagaray@gmail.com

17

Abstract

18 Human microbiota plays an important role in determining changes from health
19 to disease. Increasing research activity is dedicated to understand its diversity and
20 variability. We analyse 16S rRNA and whole genome sequencing (WGS) data from
21 the gut microbiota of 97 individuals monitored in time. Temporal fluctuations in the
22 microbiome reveal significant differences due to factors that affect the microbiota such
23 as dietary changes, antibiotic intake, early gut development or disease. Here we show
24 that a fluctuation scaling law describes the temporal variability of the system and that
25 a noise-induced phase transition is central in the route to disease. The universal law
26 distinguishes healthy from sick microbiota and quantitatively characterizes the path
27 in the phase space, which opens up its potential clinical use and, more generally, other
28 technological applications where microbiota plays an important role.

29

Importance

30 Human microbiota is tightly associated to the health status of a person.
31 Here we analyse the microbial composition of several subjects under dif-
32 ferent conditions, over a time span that ranges from days to months. Using
33 the Langevin equation as the basis of our mathematical framework in or-
34 der to evaluate microbial temporal stability, we prove that we are capable
35 to distinguish stable from unstable microbiotas. This first step will help us
36 to determine how microbiota temporal stability is related to the healthi-
37 ness of the people, and it will allow the development of a more complete
38 framework in order to deepen the knowledge of this complex system.

39 **Keywords**— microbiome, systems biology, ecological modeling, metagenomics, stability

40 **Introduction**

41 The desire to understand the factors that influence human health and cause dis-
42 eases has always been one of the major driving forces of biological research.
43 We are populated by a myriad of microorganisms that are interacting with us
44 in several physiological processes such as metabolism regulation or maturation
45 of the immune system. Human microbiota has been suggested to be closely re-
46 lated to diseases like type 2 diabetes (1), cardiovascular disease (CVD) (2),
47 irritable bowel syndrome (3), Crohn's disease (4) or some affections as obe-
48 sity (5, 6) or malnutrition (7). High throughput methods for microbial 16S
49 ribosomal RNA gene and WGS have now begun to reveal the composition of
50 archaeal, bacterial, fungal and viral communities located both, in and on the
51 human body. Modern high-throughput sequencing and bioinformatics tools
52 provide a powerful means of understanding how the human microbiome con-
53 tributes to health and its potential as a target for therapeutic interventions
54 [ref?].

55 Biology has recently acquired new technological and conceptual tools to inves-
56 tigate, model and understand living organisms at the system level, thanks to
57 the spectacular progress in quantitative techniques, large-scale measurement
58 methods and the integration of experimental and computational approaches.

59 Systems Biology has mostly been devoted to the study of well-characterized
60 model organisms but, since the early days of the Human Genome Project [ref]
61 it has become clear that applications of system-wide approaches to Human
62 Biology would bring huge opportunities in Medicine. Great effort has been
63 placed to unveil the general laws governing the behaviour of this complex sys-
64 tem [ref]. Due to his nature, microbiota can be studied under the light of the
65 ecology, where we can find general principles as the Taylor's law, which re-
66 lates spatial or temporal variability of the population with its mean. This law,
67 also known as fluctuation scale law, is ubiquitous in the natural world and can
68 be found in several systems as cosmic rays [ref1], stock markets [ref2,3], ani-
69 mal populations [ref4, 5, 6], gene expression [ref7], or in the human genome
70 [ref8]. Taylor's law has been applied to microbiota in a spatial way in the
71 work of Zhang et al., (2014), where they show that this population tend to be
72 in an aggregated way rather than in a random distribution.

73 Here we present the imprints of disease in macroscopic properties of the sys-
74 tem, by studying the temporal variability in the microbiome. We have analyzed
75 more than 35000 time series of taxa from the gut microbiome of 97 individuals
76 obtained from publicly available high throughput sequencing data on different
77 conditions: diseases, diets, obese status, antibiotic perturbation and healthy
78 individuals. Having seen that all cases follows Taylor's law, we use this em-
79 pirical fact to model how the relative abundances of taxa evolves toward time
80 thanks to the Langevin equation, in a similar way as Blumm et al., did in their
81 (2012). We use this mathematical framework to explore the temporal stability

82 of the microbiota in different conditions in order to understand how this af-
 83 fects the healthy status of the subjects. Finally, we have engineered a complete
 84 software framework, ComplexCruncher, to support the analysis of the dynam-
 85 ics of ranking processes in complex systems, which is ready to be implemented
 86 by other users.

87 **Results**

88 **Global results**

89 We have analysed the microbiome temporal variability to extract global prop-
 90 erties of the system. As fluctuations in total counts are plagued by systematic
 91 errors we worked on temporal variability of relative abundances for each taxon.
 92 Our first finding was that, in all cases, changes in relative abundances of taxa
 93 follow a ubiquitous pattern known as the fluctuation scaling law (15) or Tay-
 94 lor's power law (16), i.e., microbiota of all detected taxa follows $\sigma_i = V \cdot x_i^\beta$,
 95 a power law dependence between mean relative abundance x_i and dispersion
 96 σ_i . The law seem to be ubiquitous, spanning even to six orders of magnitude
 97 in the observed relative abundances (see Figure 1).

98 The power law (or scaling) index β and the variability V (hereafter Taylor
 99 parameters) appear to be correlated with the stability of the community and
 100 related with the health status of the host, which we consider the main finding

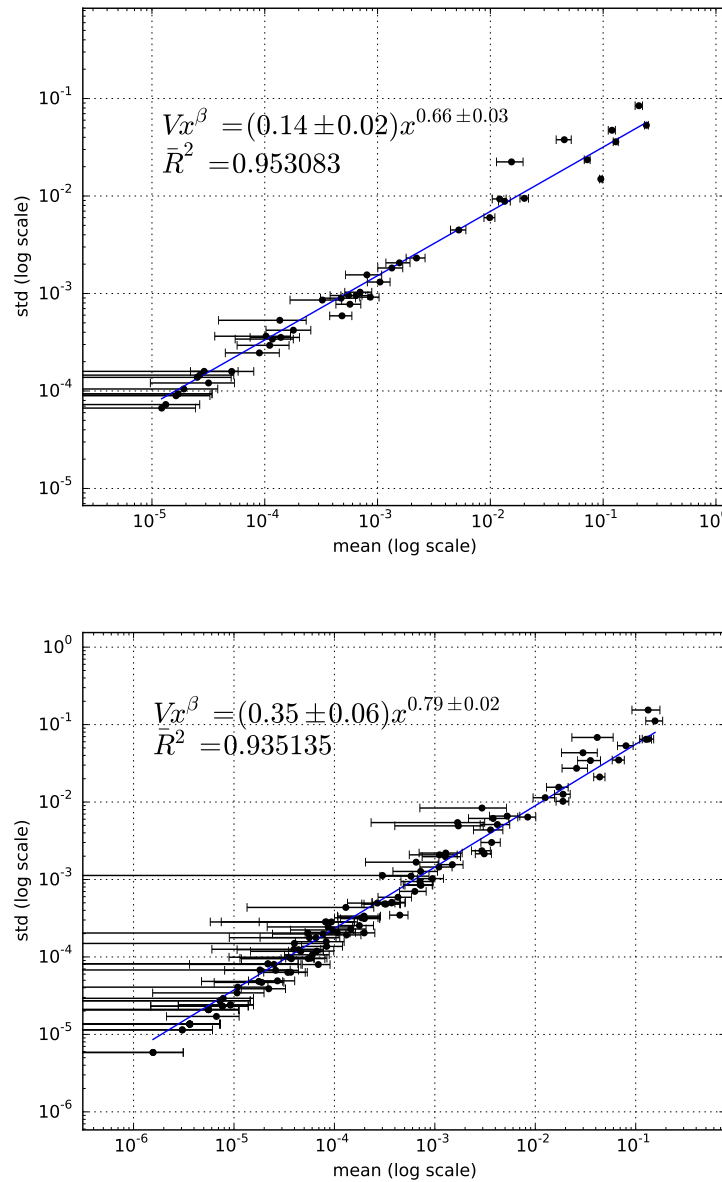


Figure 1. X-weighted power-law fits of the standard deviations versus the mean values for each bacterial genus monitored in time. We show the fit for samples from a healthy subject (top) and from a subject diagnosed with irritable bowel syndrome (bottom), studied in our lab (3). Taylor's power law seems to be ubiquitous, spanning to six orders of magnitude.

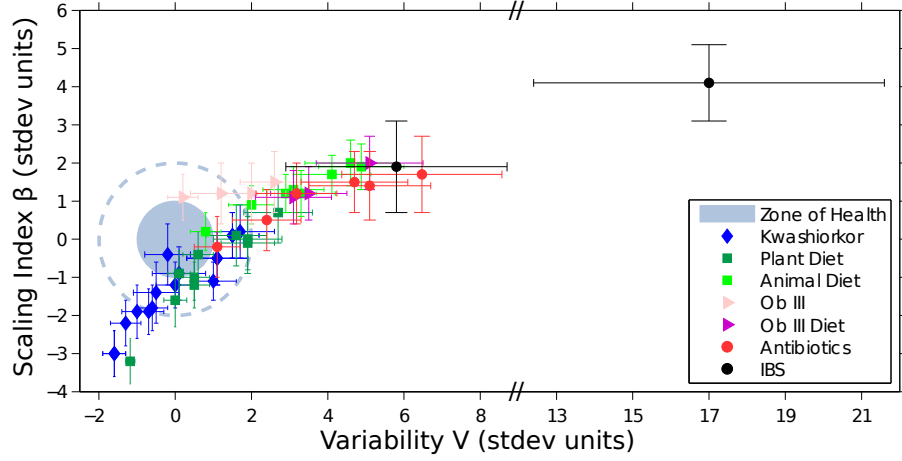


Figure 2. Taylor’s law parameter space. We have compiled here all the data studied in this work. The coloured circle corresponds to 68% confidence level (CL) region of healthy individuals in the Taylor parameter space, while dashed line delimites the 98% CL region. Points with errors place each individual gut microbiome in the Taylor space. Note that the parameters have been standardized (stdev units) to the healthy group in each study for demonstrative and comparative purposes.

101 exposed in this article (see Figure 2).

102 Taylor parameters describing the temporal variability of the gut microbiome in
 103 our sampled individuals are shown in Tables 1 to 6. Our results hint at an ubiq-
 104 uitous behaviour. On the first hand, the variability (which corresponds to the
 105 maximum amplitude of fluctuations) is large, which suggests resilient capacity
 106 of the microbiota. On the other hand, the scaling index is always smaller than
 107 one, which means that more abundant taxa are less volatile than less abundant

Metadata	V	β	\bar{R}^2	V_{st}	β_{st}
A	0.26 ± 0.05	0.826 ± 0.025	0.918	3.1 ± 0.9	1.2 ± 0.6
A	0.32 ± 0.06	0.857 ± 0.025	0.924	4.4 ± 1.1	2.0 ± 0.6
A	0.194 ± 0.033	0.813 ± 0.024	0.918	1.9 ± 0.6	0.9 ± 0.6
A	0.24 ± 0.04	0.824 ± 0.020	0.924	2.7 ± 0.7	1.2 ± 0.5
A	0.34 ± 0.06	0.855 ± 0.024	0.931	4.7 ± 1.1	1.9 ± 0.6
A	0.30 ± 0.05	0.847 ± 0.022	0.921	3.9 ± 1.0	1.7 ± 0.5
A	0.133 ± 0.021	0.784 ± 0.023	0.916	0.7 ± 0.4	0.2 ± 0.6
A	0.25 ± 0.04	0.831 ± 0.024	0.929	3.0 ± 0.8	1.4 ± 0.6
P	0.23 ± 0.05	0.804 ± 0.035	0.885	2.6 ± 0.9	0.7 ± 0.8
P	0.097 ± 0.018	0.705 ± 0.031	0.891	0.03 ± 0.34	-1.6 ± 0.7
P	0.037 ± 0.006	0.642 ± 0.025	0.881	-1.12 ± 0.11	-3.1 ± 0.6
P	0.118 ± 0.019	0.723 ± 0.025	0.895	0.4 ± 0.4	-1.2 ± 0.6
P	0.17 ± 0.04	0.78 ± 0.04	0.842	1.5 ± 0.7	0.1 ± 0.9
P	0.123 ± 0.020	0.757 ± 0.026	0.914	0.5 ± 0.4	-0.4 ± 0.6
P	0.19 ± 0.05	0.77 ± 0.04	0.871	1.8 ± 0.9	-0.0 ± 0.9
P	0.121 ± 0.020	0.736 ± 0.027	0.921	0.5 ± 0.4	-0.9 ± 0.6
P	0.187 ± 0.034	0.771 ± 0.030	0.908	1.8 ± 0.7	-0.1 ± 0.7
P	0.097 ± 0.015	0.735 ± 0.025	0.922	0.05 ± 0.28	-0.9 ± 0.6

Table 1. Taylor parameters for individuals with either animal-based (A) or plant-based (P) diets (11). Previous to diet, the population sampled is described by $\bar{V} = 0.09 \pm 0.05$, $\bar{\beta} = 0.77 \pm 0.04$, which we used to describe the *healthy zone* for this study.

ones. In addition, Taylor parameters for the microbiome of healthy individuals in different studies are compatible within estimated errors. This enables us to define an area in the Taylor parameter space that we called the *healthy zone*. In order to jointly visualize and compare the results of individuals from different studies, their Taylor parameters have been standardized, where standardization means that each parameter is subtracted by the mean value and

Metadata	V	β	\bar{R}^2	V_{st}	β_{st}
Ab	0.35 ± 0.07	0.81 ± 0.04	0.925	4.3 ± 1.4	1.3 ± 0.9
Ab	0.41 ± 0.09	0.82 ± 0.04	0.908	5.6 ± 1.8	1.6 ± 0.9
Ab	0.23 ± 0.04	0.770 ± 0.031	0.920	2.1 ± 0.8	0.5 ± 0.7
Ab	0.165 ± 0.029	0.738 ± 0.031	0.928	0.9 ± 0.6	-0.3 ± 0.7
Ab	0.34 ± 0.06	0.812 ± 0.032	0.936	4.1 ± 1.2	1.5 ± 0.7
Ab	0.26 ± 0.05	0.798 ± 0.033	0.931	2.8 ± 0.9	1.1 ± 0.8

Table 2. Taylor parameters for individuals taking antibiotics (12). Prior to antibiotics intake, the population sampled is described by $\bar{V} = 0.12 \pm 0.05$, $\bar{\beta} = 0.75 \pm 0.04$, which characterize the *healthy zone* for this study.

Metadata	V	β	\bar{R}^2	V_{st}	β_{st}
IBS	0.204 ± 0.034	0.739 ± 0.029	0.916	7.6 ± 3.7	1.9 ± 1.2
IBS	0.35 ± 0.05	0.793 ± 0.023	0.935	23.1 ± 5.9	4.0 ± 0.9

Table 3. Taylor parameters for persons diagnosed with irritable bowel syndrome (IBS) (3). Healthy individuals sampled in this study are characterized by $\bar{V} = 0.134 \pm 0.009$, $\bar{\beta} = 0.691 \pm 0.025$, which we used to define the correspondent *healthy zone*.

Metadata	V	β	\bar{R}^2	V_{st}	β_{st}
DH	0.27 ± 0.04	0.835 ± 0.016	0.925	0.2 ± 0.4	-1.0 ± 0.6
DH	0.36 ± 0.06	0.858 ± 0.015	0.929	1.1 ± 0.6	-0.2 ± 0.5
DH	0.35 ± 0.06	0.859 ± 0.014	0.926	1.0 ± 0.5	-0.1 ± 0.5
DH	0.25 ± 0.04	0.829 ± 0.014	0.911	0.0 ± 0.4	-1.2 ± 0.5
DH	0.30 ± 0.05	0.844 ± 0.014	0.920	0.5 ± 0.4	-0.7 ± 0.5
DH	0.29 ± 0.05	0.850 ± 0.016	0.915	0.4 ± 0.5	-0.5 ± 0.5
DH	0.28 ± 0.05	0.848 ± 0.016	0.921	0.3 ± 0.5	-0.5 ± 0.6
DH	0.35 ± 0.07	0.861 ± 0.017	0.918	0.9 ± 0.6	-0.0 ± 0.6
DH	0.31 ± 0.04	0.833 ± 0.012	0.916	0.6 ± 0.4	-1.1 ± 0.4
DH	0.33 ± 0.05	0.843 ± 0.013	0.925	0.8 ± 0.5	-0.7 ± 0.5
DH	0.31 ± 0.05	0.852 ± 0.014	0.925	0.6 ± 0.5	-0.4 ± 0.5
DH	0.31 ± 0.05	0.853 ± 0.015	0.930	0.6 ± 0.5	-0.4 ± 0.5
DH	0.203 ± 0.033	0.815 ± 0.015	0.907	-0.44 ± 0.32	-1.7 ± 0.5

Table 4. Taylor parameters for the healthy subject of the discordant twins (10).

This table continues in Table 5. The population of healthy twins is characterized by $\bar{V} = 0.25 \pm 0.10$, $\bar{\beta} = 0.863 \pm 0.028$, values which we used to describe the *healthy zone* for this study.

Metadata	V	β	\bar{R}^2	V_{st}	β_{st}
DK	0.40 ± 0.07	0.859 ± 0.017	0.926	1.5 ± 0.7	-0.1 ± 0.6
DK	0.44 ± 0.08	0.868 ± 0.016	0.919	1.8 ± 0.8	0.2 ± 0.6
DK	0.196 ± 0.031	0.819 ± 0.014	0.916	-0.50 ± 0.30	-1.5 ± 0.5
DK	0.160 ± 0.026	0.798 ± 0.015	0.904	-0.85 ± 0.25	-2.3 ± 0.5
DK	0.30 ± 0.05	0.845 ± 0.014	0.924	0.5 ± 0.4	-0.6 ± 0.5
DK	0.23 ± 0.04	0.834 ± 0.014	0.908	-0.1 ± 0.4	-1.0 ± 0.5
DK	0.27 ± 0.05	0.848 ± 0.015	0.930	0.2 ± 0.4	-0.5 ± 0.5
DK	0.35 ± 0.07	0.860 ± 0.019	0.916	1.0 ± 0.7	-0.1 ± 0.7
DK	0.34 ± 0.05	0.835 ± 0.012	0.917	0.9 ± 0.5	-1.0 ± 0.4
DK	0.25 ± 0.04	0.831 ± 0.012	0.912	0.0 ± 0.4	-1.1 ± 0.4
DK	0.36 ± 0.06	0.858 ± 0.013	0.918	1.1 ± 0.5	-0.2 ± 0.5
DK	0.31 ± 0.06	0.851 ± 0.016	0.924	0.6 ± 0.6	-0.4 ± 0.6
DK	0.149 ± 0.022	0.799 ± 0.013	0.905	-0.96 ± 0.22	-2.2 ± 0.5

Table 5. Taylor parameters for the kwashiorkor part of the discordant twins (10). This is a continuation of Table 4, so that the population of healthy twins is also characterized by $\bar{V} = 0.25 \pm 0.10$ and $\bar{\beta} = 0.863 \pm 0.028$.

Metadata	V	β	\bar{R}^2	V_{st}	β_{st}
OW	0.59 ± 0.12	0.894 ± 0.034	0.920	6.6 ± 2.0	2.6 ± 1.0
OW	0.22 ± 0.04	0.830 ± 0.030	0.904	0.5 ± 0.6	0.7 ± 0.9
OBI	0.28 ± 0.04	0.855 ± 0.022	0.958	1.5 ± 0.6	1.4 ± 0.6
OBI	0.33 ± 0.07	0.870 ± 0.031	0.916	2.4 ± 1.1	1.9 ± 0.9
OBII	0.223 ± 0.032	0.823 ± 0.023	0.938	0.6 ± 0.5	0.5 ± 0.7
OBII	0.208 ± 0.029	0.844 ± 0.022	0.935	0.4 ± 0.5	1.1 ± 0.7
OBIII	0.34 ± 0.05	0.855 ± 0.025	0.943	2.5 ± 0.9	1.4 ± 0.7
OBIII	0.26 ± 0.04	0.845 ± 0.026	0.954	1.1 ± 0.7	1.2 ± 0.8
OBIII	0.33 ± 0.06	0.870 ± 0.027	0.908	2.4 ± 1.0	1.9 ± 0.8
OBIII	0.200 ± 0.026	0.843 ± 0.020	0.949	0.2 ± 0.4	1.1 ± 0.6
OBIII	0.30 ± 0.05	0.846 ± 0.026	0.929	1.9 ± 0.8	1.2 ± 0.7
OBIII	0.176 ± 0.029	0.826 ± 0.026	0.894	-0.2 ± 0.5	0.6 ± 0.8
OBIII	0.30 ± 0.06	0.841 ± 0.031	0.896	1.8 ± 0.9	1.0 ± 0.9
OBIII	0.28 ± 0.04	0.857 ± 0.025	0.941	1.5 ± 0.7	1.5 ± 0.7
OBIII	0.122 ± 0.018	0.822 ± 0.024	0.930	-1.05 ± 0.30	0.5 ± 0.7
OBIIId	0.47 ± 0.08	0.872 ± 0.023	0.945	4.7 ± 1.3	1.9 ± 0.7
OBIIId	0.38 ± 0.06	0.846 ± 0.023	0.951	3.2 ± 1.0	1.2 ± 0.7
OBIIId	0.36 ± 0.06	0.842 ± 0.022	0.954	2.9 ± 0.9	1.1 ± 0.6

Table 6. Taylor parameters for individuals with different degrees of overweight and obesity (9). Healthy people in this study, whom were not obese, are characterized by $\bar{V} = 0.19 \pm 0.06$, $\bar{\beta} = 0.806 \pm 0.034$, which we used to determine the correspondent *healthy zone* for this study.

divided by the standard deviation of the group of healthy individuals for each study (for details of the procedure, please see Standardization subsection in Material and Methods). The healthy zone and the standardized Taylor parameters for individuals whose gut microbiota is threatened (i.e., suffering from kwashiorkor, altered diet, antibiotics or IBS) is shown in Figure 2. Children developing kwashiorkor show smaller variability than their healthy twins. A meat/fish-based diet increases the variability significantly when compared to a plant-based diet. All other cases presented increased variability, which is particularly severe, and statistically significant at more than 95% CL, for obese patients grade III on a diet, individuals taking antibiotics or IBS–diagnosed patients. A global property emerges from all worldwide data collected: Taylor parameters characterize the statistical behaviour of microbiome changes. Furthermore, we have verified that our conclusions are robust to systematic errors due to taxonomic assignment.

Taylor’s power law has been explained in terms of various effects, all without general consensus. It can be shown to have its origin in a mathematical convergence similar to the central limit theorem, so virtually any statistical model designed to produce a Taylor law converge to a Tweedie distribution (17), providing a mechanistic explanation based on the statistical theory of errors (18–20). To unveil the generic mechanisms that drive different scenarios in the β – V space, we model the system by assuming that taxon relative abundance follows a Langevin equation with, on the one hand, a deterministic term that captures the fitness of each taxon and, on the other hand, a ran-

domness term associated with Gaussian random noise (21). Both terms are modeled by power laws, with coefficients that can be interpreted as the taxon fitness F_i and the variability V (see Model under Material and Methods). In this model, when V is sufficiently low, abundances are stable in time. Differences in variability V can induce a noise-induced phase transition in relative abundances of taxa. The temporal evolution of the probability of a taxon having abundance x_i given its fitness is governed by the Fokker–Planck equation. The results of solving this equation show that the stability is best captured by a phase space determined by fitness F and amplitude of fluctuations V (see Figure 3).

The model predicts two phases for the gut microbiome: a stable phase with large variability that permits some changes in the relative abundances of taxa and an unstable phase with larger variability, above the phase transition, where the order of abundant taxa varies significantly with time. The microbiome of all healthy individuals was found to be in the stable phase, while the microbiome of several other individuals was shown to be in the unstable phase. In particular, individuals taking antibiotics and IBS–diagnosed patient P2 had the most severe symptoms. In this phase diagram, each microbiota state is represented by a point at its measured variability V and inferred fitness F . The model predicts high average fitness for all taxa, i.e., taxa are narrowly distributed in F . The fitness parameter has been chosen with different values for demonstrative purposes. Fitness is larger for the healthiest subjects and smaller for the IBS–diagnosed patients.

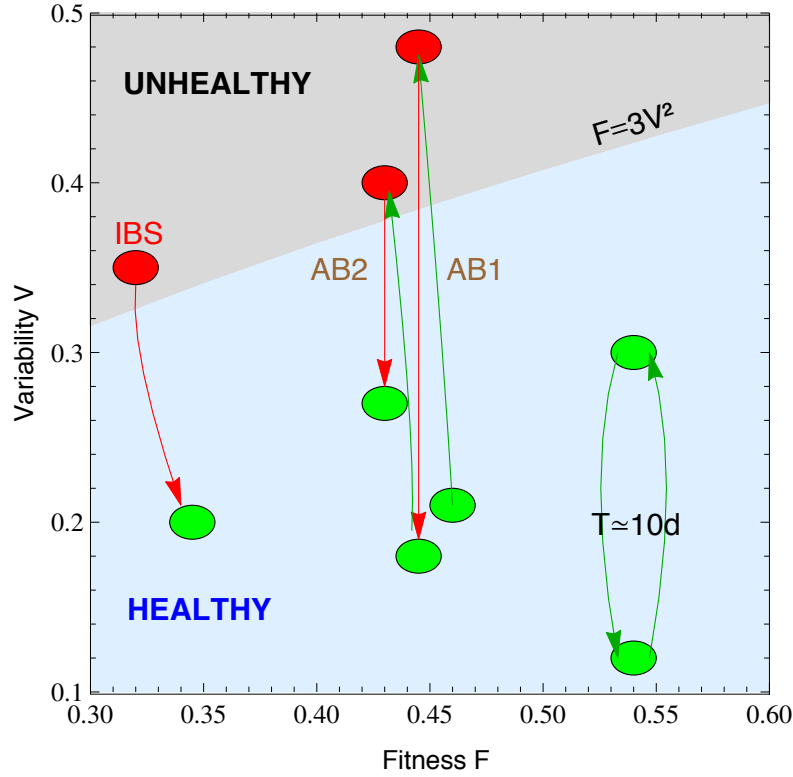


Figure 3. Microbiota states can be placed in the phase space F - V . The light blue shaded region corresponds to the stable phase, while the grey shaded region is the unstable phase (the phase transition line is calculated for $\alpha = \beta = 0.75$). We place healthy individuals (green) and individuals whose gut microbiota is threatened (antibiotics, IBS) in the phase space fitness–variability. Gut microbiota of healthy individuals over a long term span show a quasi-periodical variability (central period is ten days). We show that taking antibiotics (AB1 and AB2 correspond to first and second treatment respectively) induces a phase transition in the gut microbiota, which impacts its future changes. We also show an IBS–diagnosed patient transiting from the unstable to the stable phase.

160 **Specific results**

161 **Fit Plots**

162 For each and every dataset included in the study, an unweighted fit (see Mate-
 163 rial and Methods section for details) and a X-weighted fit (detailed in Material
 164 and Methods subsection) have been calculated for standard deviations versus
 165 the mean values for each bacterial genus monitored in time. Figure 1 showed
 166 the X-weighted fit for samples from a healthy subject (top) and from a sub-
 167 ject diagnosed with irritable bowel syndrome (bottom) studied in our lab (3),
 168 while Figure 4 shows the corresponding unweighted fits. Additionally, for the
 169 unweighted fit, a complete residues analysis is performed, and a 4-in-1 figure
 170 is generated as shown in Figure 5, corresponding to patient A (top plot in Fig-
 171 ure 1 and Figure 4). Among other tests, it allows to check for normality and
 172 homoscedasticity of the residues.

173 **Histogram Plots (DRAFT)**

174 *cmplxcruncher* generates three different histogram plots:

175 —**Absolute frequencies plot** : This plot is useful to visually assess the validity
 176 of the time points in terms of the accumulated absolute frequency of the
 177 elements (taxa), since absolute frequencies far (much higher or much
 178 lower) from those typically observed could mean a sampling problem. As
 179 an example, Figure ?? shows this histogram for the pre-treatment data

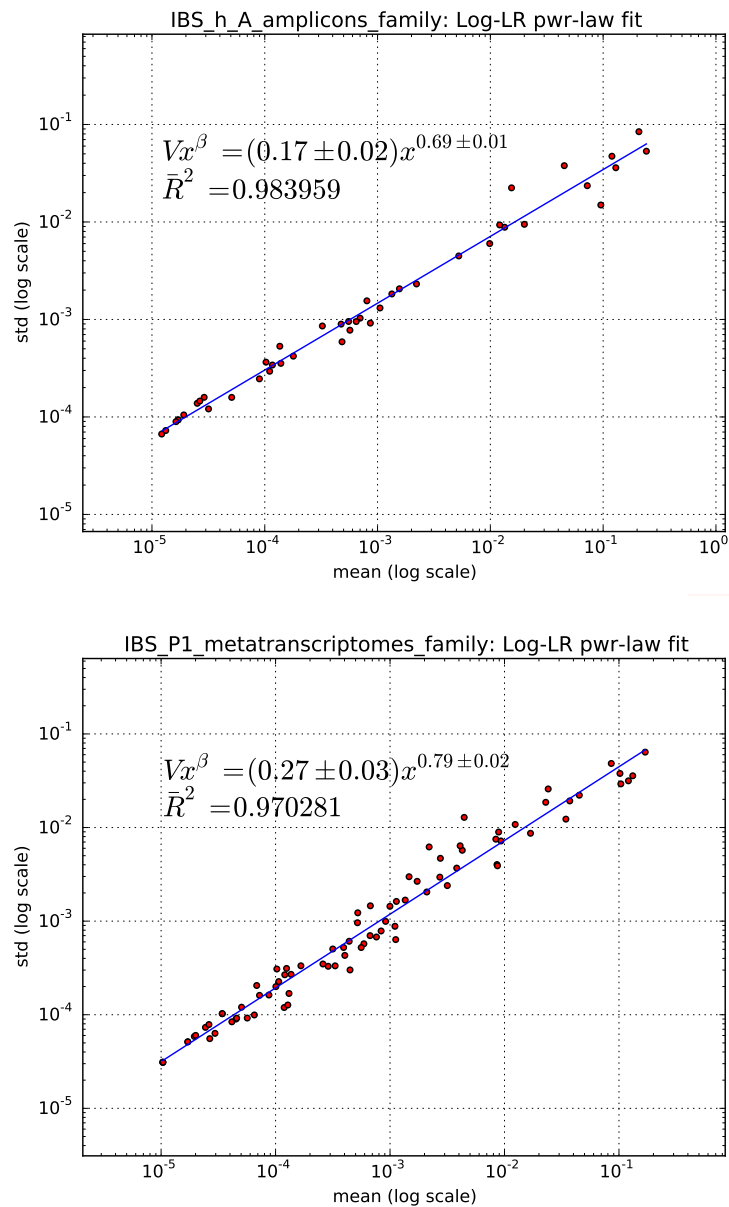


Figure 4. Log plots of unweighted fits corresponding to the datasets shown in Figure 1

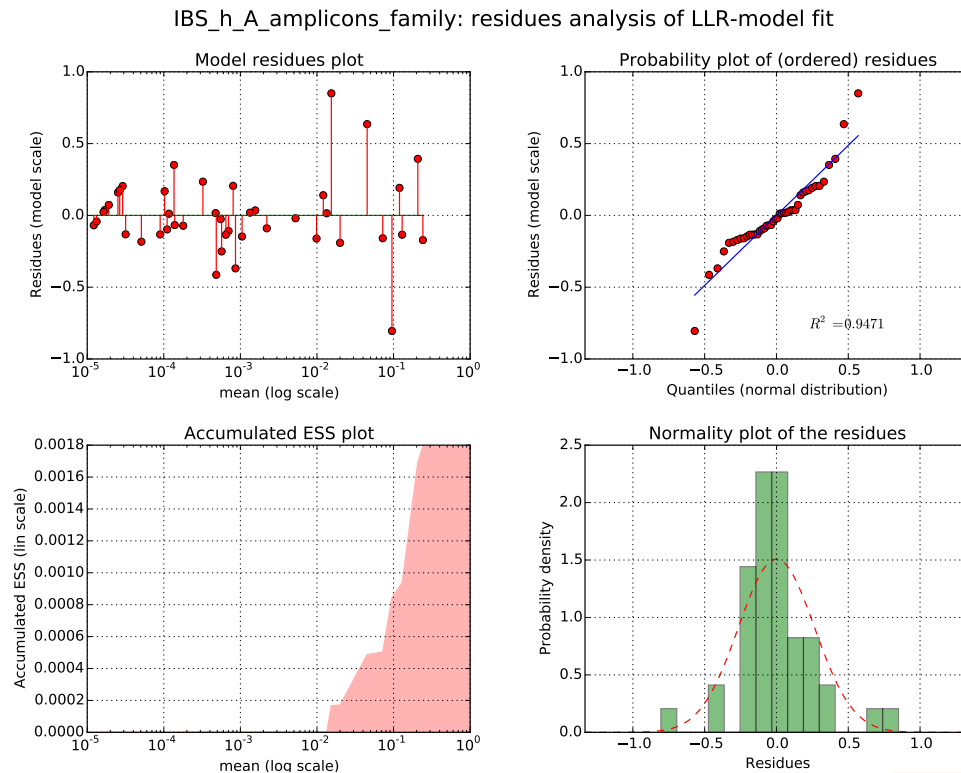


Figure 5. Residues analysis plot corresponding to the unweighted fit for patient A of the IBS study. The top-left subplot is a simple residues plot. The top-right subplot is a Normal quantiles plot with linear fitting (value of coefficient of determination is provided). The bottom-left subplot shows an accumulated ESS (Explained Sum of Squares) plot. Finally, the bottom-right subplot is a residues Normal histogram plot. This set of subplots allows to check for normality and homoscedasticity of the residues.

(first 7 times) of patient “D” in the antibiotics study (12).

—**2D deviation plot** : The 2D semi-logarithmic histogram representing deviations from the mean versus the mean itself, is a useful tool in the analysis of the stability of ranking processes in complex systems (21). Figure 6 shows this plot for the data used in the fit shown in Figure ??.

—**Zero relative frequency plot** : We could define the ZRF (Zero Relative Frequency, thereby ranging from 0 to 1) of an element (taxon) as the portion of times where it is zero, i.e., it is not found. Attending to all the elements (taxa), we can plot the ZRF histogram, which then lies on the horizontal axis of the plot. The vertical axis shows the number of elements (taxa), so the height of a bar represents the number of elements that have determinate ZRF. In this respect, the bar over 0.0 counts the number of elements (taxa) that are present at every time point of the data set (aka “core”), while the bar over 1.0 would count the number of elements (taxa) that are never found (this bar never appears because all these “null” elements are automatically filtered by the code). Figure 7 shows an example of this plot. There, we can see that 12 taxa are present at all the time points of the time series while 9 taxa basically appear only once. So, this plot is clearly useful to notice how the “core” is distributed.

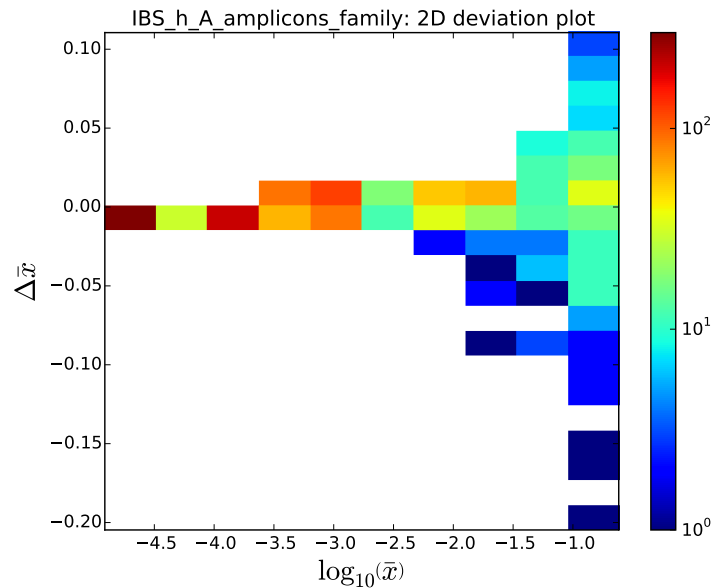


Figure 6. 2D histogram deviation plot of the data for patient A of the IBS study

199 Correlation and Rank Plots

200 *cmplxcruncher* generates three different plots falling under this category, as
 201 well as Excel files with the resulting matrices:

202 **—Elements correlation matrix** : This plot shows a correlation matrix among
 203 the elements (taxa), calculated with the time as independent variable.
 204 For these calculations, the data set is not normalized to avoid entering
 205 an additional constraint. As an example, Figure 8 shows this matrix for
 206 the “core” elements (taxa) present in the pre-treatment data (first seven
 207 times) of patient “D” in the antibiotics study (12).

208 **—Times correlation matrix** : This plot presents a correlation matrix among

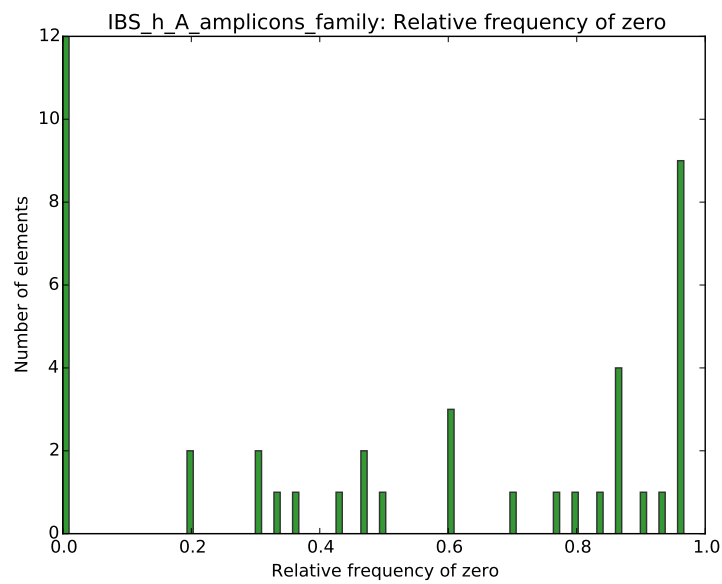


Figure 7. Histogram with the relative frequency of zero for the elements (taxa) in the data for patient A of the IBS study

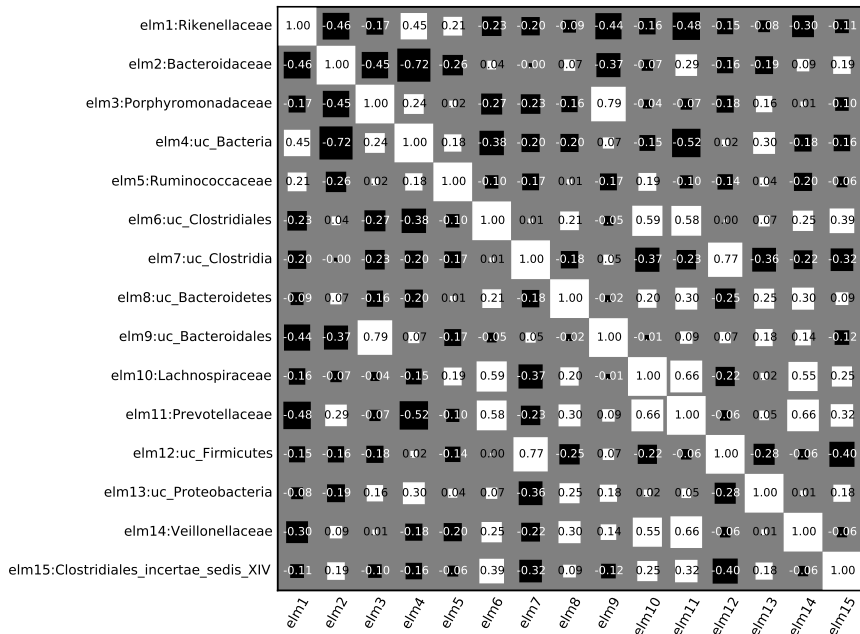


Figure 8. Element correlation plot of the for the most dominant taxa in the data for patient A of the IBS study

the time points of the data set, calculated with the elements (taxa) as independent variable. Again, the data set is not normalized. Figure ?? shows this matrix for the “core” elements (taxa) present in pre-treatment data (first seven times) of patient “D” in the antibiotics study (12).

—**Rank dynamics and stability plot** : This plot shows the variation in the rank with time for the most dominant elements (taxa) and their calculated RSI, as discussed in Section . Figure 9 shows this plot for the elements (taxa) in the data used in the fit shown in Figure ??.

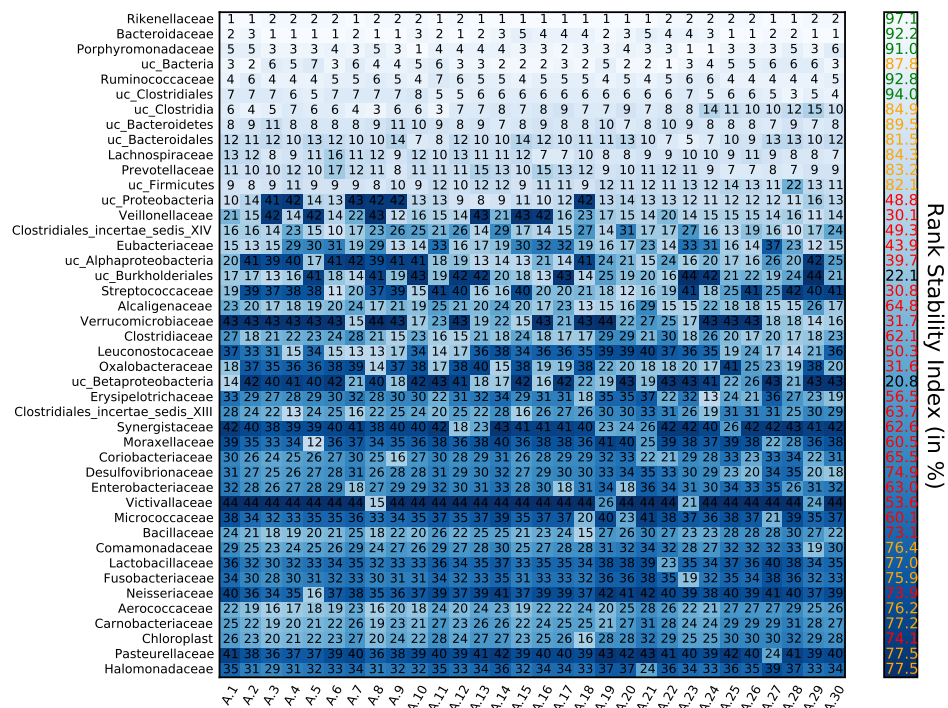


Figure 9. Matrix showing the rank variation throughout time for the most dominant elements (taxa) and their calculated Rank Stability Index (as discussed in Material and Methods) in the data used in the fit shown in Figure ??

217 Time dependence of model parameters

218 Finally, we have studied the time dependence of the variability V and power
 219 law index β (see Model under Material and Methods) by using a sliding win-
 220 dow approach. The total number of time points are divided in subsets of five
 221 points, where next subset is defined by adding next time sampling and by
 222 eliminating the earliest one. Both parameters were calculated for each sub-
 223 set against the average time lapse. Figure 10 shows the variability V as a
 224 function of time for the largest sampling: two individuals in the Caporaso's
 225 study (8) corresponding to the gut microbiota of a male (upper plot) and a
 226 female (lower plot). Figure 11 shows the time evolution of V for patient P2
 227 of the IBS study (3) (upper plot) and patient D in the antibiotics study (12)
 228 (lower plot).

229 Discussion

230 We have quantitatively characterized whether the microbiota belongs to a healthy
 231 individual or a subject corresponding to an altered or pathological state (i.e.,
 232 altered diet, antibiotic treatment, early gut development, diagnosed IBS). De-
 233 ciphering the mechanisms of disease requires in depth knowledge of the un-
 234 derlying biological mechanisms. We describe here the macroscopic behavior
 235 of disease by a noise-induced phase transition with a control parameter that
 236 can be measured by the temporal variability of the microbiome. The micro-

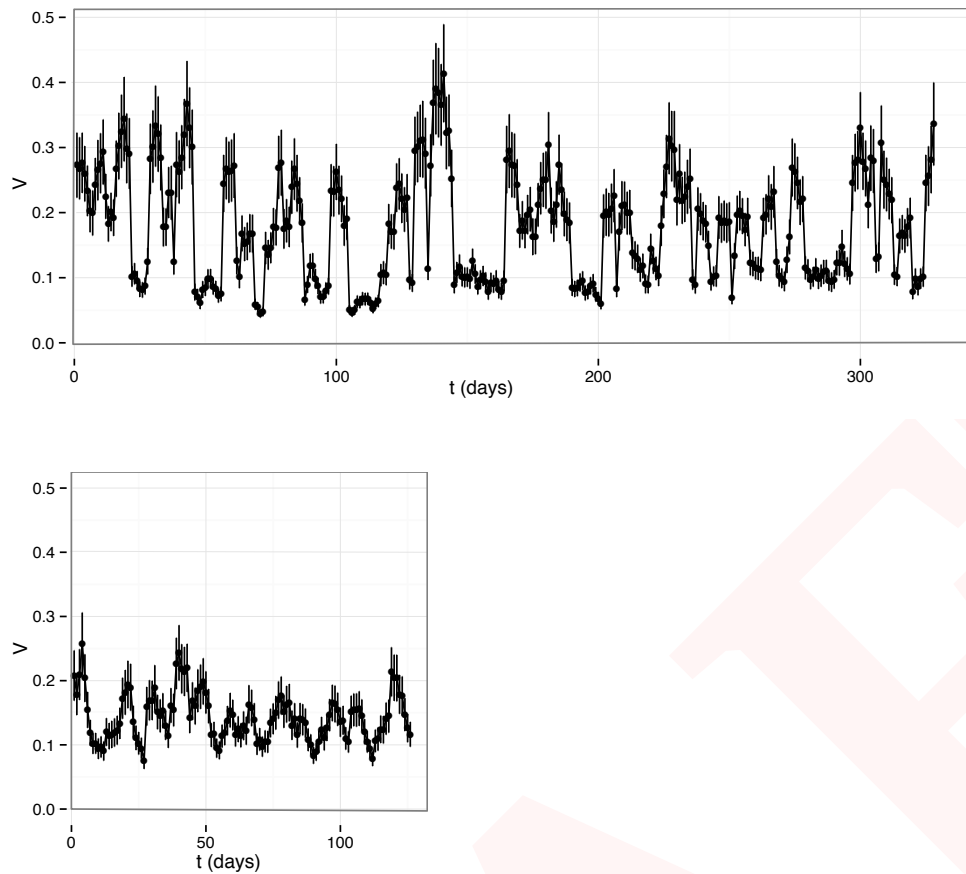


Figure 10. V as a function of time for the two individuals in the Caporaso's study (8): samples of gut microbiome of a male (upper plot) and a female (lower plot). Both samples show changes in the variability V with quasi-periodic behavior peaked at about 10 days. Variability grows more for the gut microbiota of the male and share a minimal value around 0.1 with the gut microbiota of the female.

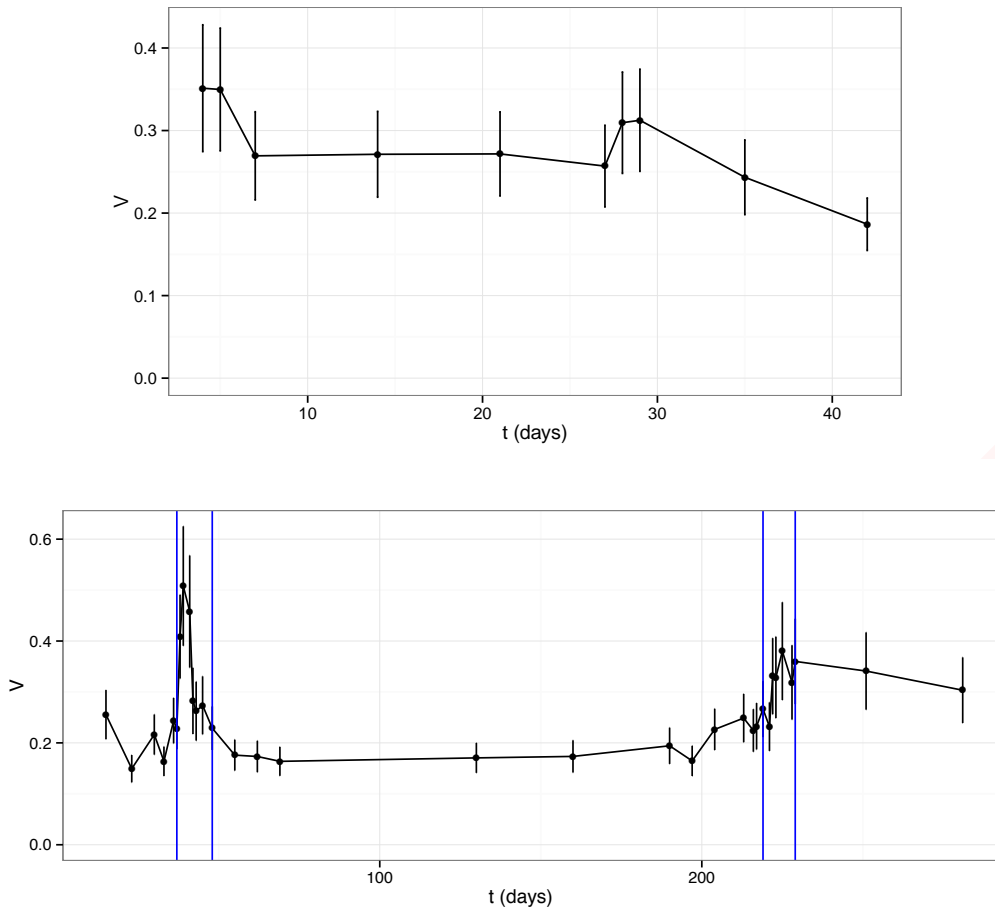


Figure 11. V as a function of time for patient P2 of the IBS study (3) (upper plot) and patient D in the antibiotics study (12) (lower plot). The variability of the gut microbiota of P2 decreases from above 0.3 to below 0.2, showing a slow tendency to increase the order of the system. Antibiotic intake leads to a quick increase of variability which lasts for a few days to recover ordering. The second antibiotic treatment shows some memory (lower increase of variability) with a slower recovery. NOTE: The blue vertical lines in the lower plot are showing the periods of antibiotic treatment.

237 biota of healthy individuals and of individuals with pathologies represent dif-
 238 ferent phases separated by this noise-induced phase transition. Improved high-
 239 throughput sequencing of samples from individuals monitored over time and
 240 taxonomic assigning methods will provide a better distinction among patholo-
 241 gies or altered states of the microbiota.

242 Materials and Methods

243 Model

244 We model the microbial abundances across time along the lines of Blumm *et al.*
 245 (21). The dynamics of taxon relative abundances is described by the Langevin
 246 equation:

$$247 \quad \dot{x}_i = F_i \cdot x_i^\alpha + V \cdot x_i^\beta \xi_i(t) - \phi(t) \cdot x_i, \quad (1)$$

248 where F_i captures the fitness of the taxon i , V corresponds to the noise am-
 249 plitude and $\xi_i(t)$ is a Gaussian random noise with zero mean $\langle \xi_i(t) \rangle = 0$
 250 and variance uncorrelated in time, $\langle \xi_i(t) \xi_i(t') \rangle = \delta(t' - t)$, . The function
 251 $\phi(t)$ ensures the normalization at all times, $\sum x_i(t) = 1$, and corresponds to
 252 $\phi(t) = \sum F_i x_i^\alpha + \sum V x_i^\beta \xi_i(t)$. The temporal evolution of the probability that
 253 a taxon i has a relative abundance $x_i(t)$, $P(x_i, t)$, is determined by the Fokker-

254 Planck equation:

$$255 \quad \frac{\partial P}{\partial t} = -\frac{\partial}{\partial x_i} [(F_i \cdot x_i^\alpha - \phi(t) \cdot x_i) \cdot P] + \frac{1}{2} \frac{\partial^2}{\partial x_i^2} (V^2 \cdot x_i^{2\beta} \cdot P). \quad (2)$$

256 The microbiota evolves towards a steady-state with a time-independent proba-
 257 bility depending on the values of α , β , F_i and V . For $\alpha < 1$ (otherwise, systems
 258 are always unstable), the steady-state probability may be localized in a region
 259 around a preferred value or broadly distributed over a wide range, depending
 260 on whether the fitness F_i dominates or is overwhelmed by the noise amplitude
 261 V . The steady-state solution of the Fokker-Planck equation is given by:

$$262 \quad P_0(x_i) = C_{ne}(\alpha, \beta, F_i, V) \cdot x_i^{-2\beta} \cdot \exp\left[\frac{2F_i}{V^2} \frac{x_i^{1+\alpha-2\beta}}{1+\alpha-2\beta} - \frac{\phi_0}{V^2} \frac{x_i^{2-2\beta}}{1-\beta}\right] \quad \text{if } 2\beta \neq 1+\alpha,$$

$$263 \quad P_0(x_i) = C_e(\alpha, \beta, F_i, V) \cdot x_i^{\frac{2F_i}{V^2}-2\beta} \cdot \exp\left[\frac{\phi_0}{V^2} \frac{x_i^{2-2\beta}}{1-\beta}\right] \quad \text{if } 2\beta = 1+\alpha,$$

264 where $\phi_0 = (\sum_i F_i^{1/(1-\alpha)})^{1-\alpha}$ and C_{ne} and C_e are integrals that should be solved
 265 numerically for the parameters of interest. The ordered phase happens when
 266 the solution has a maximum in the physical interval ($0 < x_i < 1$). For larger V ,
 267 the transition to a disordered phase happens when the maximum shifts to the
 268 unphysical region $x_i < 0$, which sets the phase transition region $V(\alpha, \beta, F_i)$.
 269 The phase transition region can be calculated analytically in particular cases:

$$270 \quad F_i^2 = 4\beta\phi_0V^2 \quad \text{if } \beta = \alpha \neq 1,$$

$$271 \quad F_i = \beta V^2 \quad \text{if } 2\beta = 1 + \alpha,$$

where the first case, simplifies to $F = 3V^2$ if $\beta = 0.75$ and the fitness of this taxon dominates in ϕ_0 . In many physical systems (Brownian motion is the classical example), the two terms of the Langevin equation are related. The *fluctuation–dissipation theorem* states a general relationship between the response to an external disturbance and the internal fluctuations of the system (22). The theorem can be used as the basic formula to derive the fitness from the analysis of fluctuations of the microbiota, assuming that it is in equilibrium (the ordered phase).

Explain better the fluctuation-dissipation theorem

Selection and Methods

The bacteria and archaea taxonomic assignments were obtained by analysing 16S rRNA sequences, which were clustered into operational taxonomic units (OTUs) sharing 97 % sequence identity using QIIME (13). WGS data (10) were analysed and assigned at strain level by the Livermore Metagenomic Analysis Toolkit (LMAT) (14), according to their default quality threshold. Genus, with best balance between error assignment and number of taxa, was chosen as our reference taxonomic level. We have verified that our conclusions are not significantly affected by selecting family or species as the reference taxonomic level (see Figure 12).

Specify, in each study treated, the nature of the samples (conditions,

timespan between timepoints, subjects). Specify, and it is very important, what we consider *healthy* in each study (for example: pre-antibiotics is healthy)

Sample selection

We have chosen studies about relevant pathologies containing metagenomic sequencing time data series of bacterial populations from humans in different healthy and non-healthy states. We have selected only those individuals who had three or more time points of data available in databases. Metadata of each study is provided in Tables 1 to 6. All used 16S rRNA gene sequencing except for the study of the discordant kwashiorkor twins (10) (see Tables 4 and 5) where shotgun metagenomic sequencing (SMS) and 16S rRNA were used. In the latter case we selected to work with SMS data to show that our method is valid regardless of the source of taxonomic information. Each one of the datasets was treated as follows:

16rRNA sequences processing

Reads from the selected studies were first quality filtered using the FastX toolkit (23), allowing only those reads which had more than 25 of quality along the 75% of the complete sequence. 16S rRNA reads were then clustered at 97% nucleotide sequence identity (97% ID) into operational taxonomic units (OTUs) using QIIME package software (13) (version 1.8) We followed open reference

OTU picking workflow in all cases. The clustering method used was uclust, and the OTUs were matched against Silva database (24) (version 111, July 2012) and were assigned to taxonomy with an uclust-based consensus taxonomy assigner. The parameters used in this step were: similarity 0.97, prefilter percent id 0.6, max accepts 20, max rejects 500.

Metagenomic sequences processing

Metagenomic shotgun (and 16S too) sequences were analyzed with LMAT (Livermore Metagenomics Analysis Toolkit) software package (14) (version 1.2.4, with Feb'15 release of data base *LMAT-Grand*). LMAT was run using a Bull shared-memory node belonging to the team's HPC (high performance computing) cluster. It is equipped with 32 cores (64 threads available using Intel Hyper-threading technology) as it has 2 Haswell-based Xeons, the E5-2698v3@2.3 GHz, sharing half a tebibyte (0.5 TiB, that is, 512 gibibytes) of DRAM memory. This node is also provided with a card PCIe SSD as NVRAM, the P420m HHHL, with 1.4 TB, and 750000 reading IOPS, 4 KB, achieving 3.3 GB/s, which Micron kindly issued free of charge, as a sample for testing purposes. The computing node was supplied with a RAID-0 (striping) scratch disk area. We used the "Grand" database¹, release Feb'15, provided by the LMAT team. Previously to any calculation, the full database was loaded in

¹In this context, "Grand" refers to a huge database that contains k-mers from all viral, prokaryote, fungal and protist genomes present in the NCBI database, plus Human reference genome (hg19), plus GenBank Human, plus the 1000 Human Genomes Project (HGP). This represent about 31.75 billion k-mers occupying 457.62 GB.

the NVRAM. With this configuration the observed LMAT sustained sequence classification rate was 20 kpb/s/core. Finally, it is worth mentioning that a complete set of Python scripts have been developed as back-end and front-end of the LMAT pipeline in order to manage the added complexity of time series analysis.

Taxa level selection

We selected genus as taxonomic level for the subsequent steps of our work. In order to ensure that, between adjacent taxonomic levels, there were not crucial differences which could still be of relevance after standardization (see Section), we tested two different data sets. In the former, the antibiotics study (12) with 16S data, we tested the differences between genus and family levels. The latter dataset tested was the kwashiorkor discordant twins study (10) for both genus and species taxonomic levels. The Figures 12 (overview) and 13 (detail) plot the comparison between studies (and so, 16S and SMS) and between adjacent taxonomic levels.

ComplexCruncher

A complete software framework, named 'ComplexCruncher', has been engineered to support the analysis of the dynamics of ranking processes in complex systems. Although the software was devised with a clear bias towards metage-

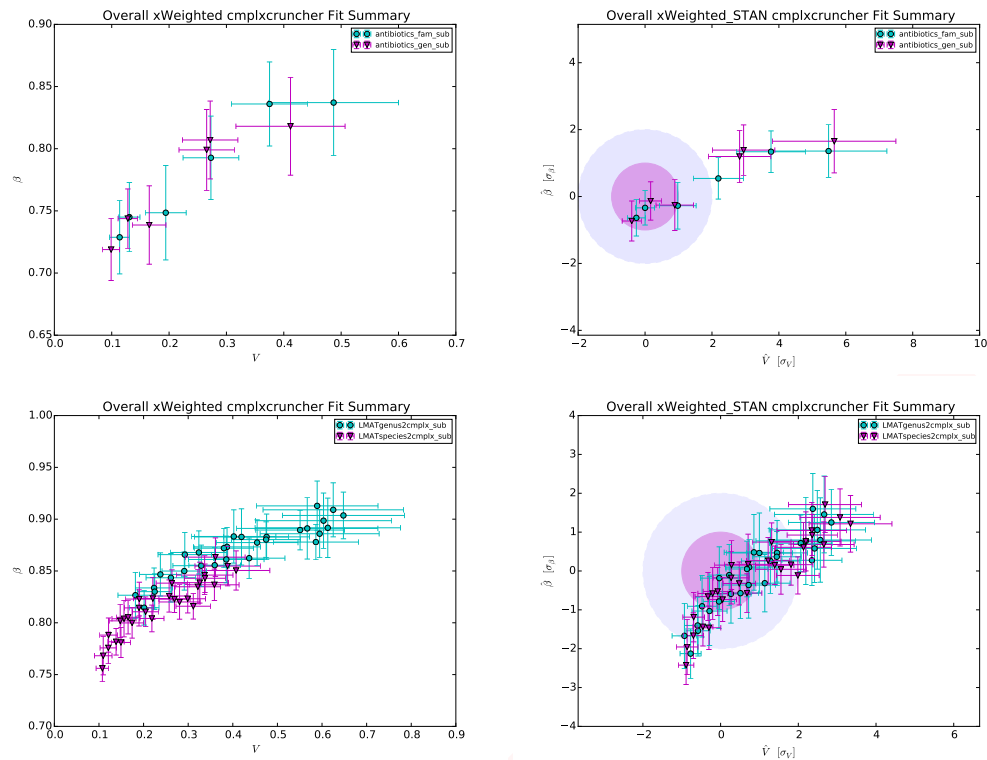


Figure 12. Overview of comparison of different approaches based on adjacent taxonomic levels using plots in the Taylor-parameters space. For 16S (former row of subfigures), the levels are family vs. genus, whereas for SMS (latter row of subfigures) levels are genus vs. species. The left column shows the raw results and the right column plots the standardized results (see Section)

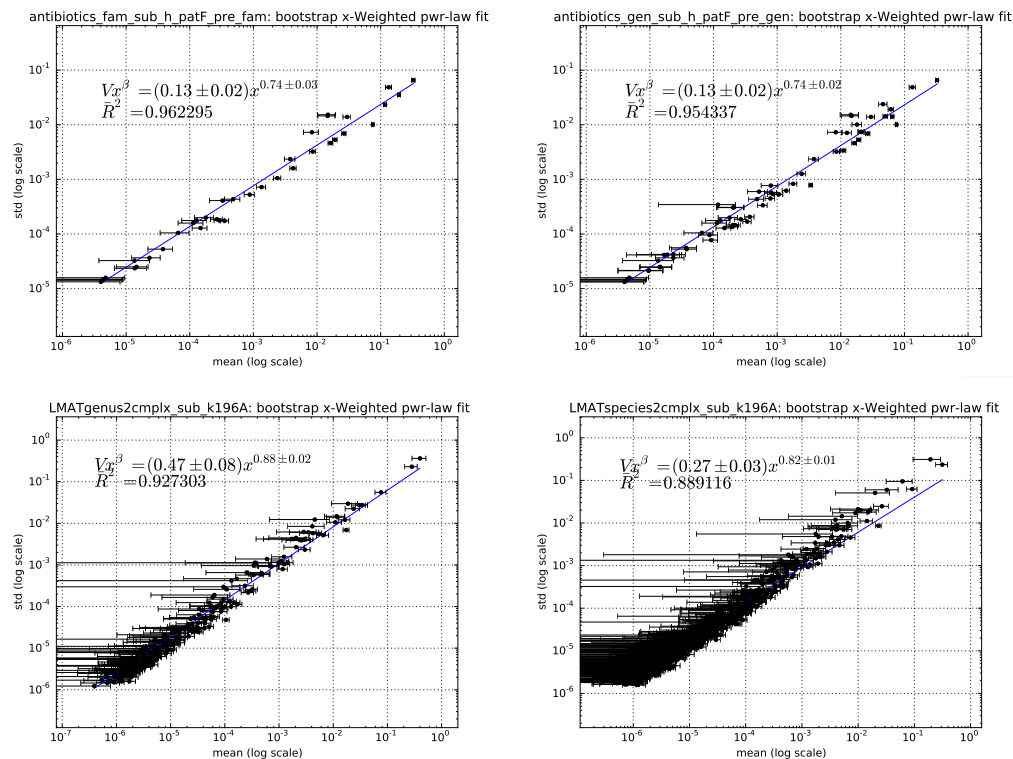


Figure 13. Detail of comparison of different approaches based on adjacent taxonomic levels using plots of X-weighted power-law fits (see Material and Methods). The former row of subfigures shows examples for 16S, whereas the latter row of subfigures plots examples for SMS. The left column shows results for the superior taxonomic level (family for 16S, genus for SMS), while the right column shows results for the inferior level (genus for 16S, specie for SMS).

nomics, it is general enough to be able to cope with a ranking process in any complex system. Implemented in Python using well-known open-source community software, the software solution is composed of two parts that can be used together or apart: a web-based graphic front-end connected to a database, and a computing kernel. Used together, this software enables other users to reproduce our results easily and, furthermore, upload and analyse their own data or experiment with the preloaded metagenomics data sets.

‘ComplexCruncher WebPortal’ (CCWebPortal) is a web platform designed to allow the user to interact with a data repository of selected and well-documented metagenomics data sources. Through a few simple steps, the user can perform advanced searches on the complete set of records in the metagenomics repository. The web application provides advanced filters that allow the user to reduce the search to a small set of interest. After this first step, the user can refine the search and discard those records that do not meet certain requirements.

The web application allows calculations to be done directly by the stable release of the *cmplxcruncher* computing kernel. At the end of the calculations, the results are displayed to the user on the same browser which runs the web application. Then, the user can interact over the series of generated graphics thus allowing flexible comparison among them. In addition, CCWebPortal enables direct download of generated data (plots, spreadsheets, etc). The web application generates a report file summarizing all the results in PDF format. If the user has login permissions, CCWebPortal enables the option of insert new

372 database records in addition to editing and deleting existing ones.

373 CCWebPortal is a web application that runs on current versions of many browsers.

374 Additional software is not needed and only requires javaScript to be enabled

375 on the browser to run applications. CCWebPortal is implemented following the

376 client–server distributed programming model, where the javaScript client ap-

377 plication connects to a remote server that enables the execution of calculations

378 and transactions through a centralized database management system. A set of

379 relational tables allows the structuring of the metagenomics repository to es-

380 tablish relationships between records. Thus the search and information thresh-

381 ing is optimized for queries launched from the client interface. Access to the

382 database on the server is implemented through Django framework, an open-

383 source framework written in Python using the model-view-controller (MVC)

384 architectural pattern for implementing user interfaces.

385 The effective data analysis has been performed with a Python tool developed

386 from scratch to more than 4200 lines of code. Implemented following the

387 Object Oriented Programming paradigm, this software is the back-end of the

388 website described above. However, it could be run as an independent piece of

389 software since it is built as a Python package provided with a command-line

390 front-end (*cmplxcruncher.py*). Once installed, the tool can be run interactively

391 but also in automatic mode, which uses parallel computation to speed up the

392 analysis of several data sources.

393 *cmplxcruncher* performs the power-law fit described in the *Blumm, N. et al.*

paper, but by fitting the best model, i.e. choosing between fitting a power-law using linear regression versus nonlinear regression (25). In the power-law fit plots we also show the generalized coefficient of determination computed for continuous models (26, 27).

Un-weighted power-law fit

Fitting the best model

As already mentioned, to choose between fitting power laws ($y = Vx^\beta$) using linear regression on log-transformed (LLR) data versus non-linear regression (NLR), we mainly follow *General Guidelines for the Analysis of Biological Power Laws* (25). It consists of the following three steps:

1. Determining the appropriate error structure by likelihood analysis.
 - (a) Fit the Non-Linear Regression (NLR) model and obtain V_{NLR} , β_{NLR} and σ_{NLR}^2 .
 - (b) Calculate the loglikelihood that the data (n is sample size) are generated from a normal distribution with additive error:
 - The likelihood of a normal distribution is:

$$\mathcal{L}_{\text{norm}} = \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma_{\text{NLR}}^2}} \exp \left(-\frac{(y_i - V_{\text{NLR}}x_i^{\beta_{\text{NLR}}})^2}{2\sigma_{\text{NLR}}^2} \right) \right]$$

- 411 • So, the loglikelihood of a normal distribution is:

$$\begin{aligned}
 412 \quad \log \mathcal{L}_{\text{norm}} &= -\frac{n}{2} \log |2\pi\sigma_{\text{NLR}}^2| - \frac{1}{2\sigma_{\text{NLR}}^2} \underbrace{\sum_{i=1}^n (y_i - V_{\text{NLR}} x_i^{\beta_{\text{NLR}}})^2}_{\text{RSS}_{\text{NLR}}} \\
 413 \quad &= -\frac{n}{2} \log |2\pi\sigma_{\text{NLR}}^2| - \frac{\text{RSS}_{\text{NLR}}}{2\sigma_{\text{NLR}}^2}
 \end{aligned}$$

- 414 (c) Calculate the *corrected Akaike's Information Criterion* for the NLR model:

$$415 \quad \text{AIC}_{\text{cNLR}} = 2k - 2 \log \mathcal{L}_{\text{norm}} + \frac{2k(k+1)}{n-k-1}$$

- 416 (d) Fit the Log-transformed Linear Regression (LLR) model and obtain

$$417 \quad V_{\text{LLR}}, \beta_{\text{LLR}} \text{ and } \sigma_{\text{LLR}}^2.$$

- 418 (e) Calculate the loglikelihood that the data (n is sample size) are gen-
 419 erated from a lognormal distribution with multiplicative error:

- 420 • The likelihood of a lognormal distribution is:

$$421 \quad \mathcal{L}_{\text{logn}} = \prod_{i=1}^n \left[\frac{1}{y_i \sqrt{2\pi\sigma_{\text{LLR}}^2}} \exp \left(-\frac{(\log |y_i| - \log |V_{\text{LLR}}| - \beta_{\text{LLR}} \log |x_i|)^2}{2\sigma_{\text{LLR}}^2} \right) \right]$$

- 422 • So, the loglikelihood of a lognormal distribution is:

$$\begin{aligned}
 \log \mathcal{L}_{\log n} &= -\frac{n}{2} \log |2\pi\sigma_{\text{LLR}}^2| - \sum_{i=1}^n \log |y_i| - \\
 &\quad - \frac{1}{2\sigma_{\text{LLR}}^2} \underbrace{\sum_{i=1}^n (\log |y_i| - \log |V_{\text{LLR}}| - \beta_{\text{LLR}} \log |x_i|)^2}_{\text{RSS}_{\text{LLR}}} \\
 &= -\frac{n}{2} \log |2\pi\sigma_{\text{LLR}}^2| - \frac{\text{RSS}_{\text{LLR}}}{2\sigma_{\text{LLR}}^2} - \sum_{i=1}^n \log |y_i|
 \end{aligned}$$

- 426 (f) Calculate the *corrected Akaike's Information Criterion* for the LR model:

$$\text{AIC}_{\text{cLLR}} = 2k - 2 \log \mathcal{L}_{\log n} + \frac{2k(k+1)}{n-k-1}$$

428 2. Compare AIC_{cNLR} with AIC_{cLLR} :

- 429 • If $\text{AIC}_{\text{cNLR}} - \text{AIC}_{\text{cLLR}} < -2$, the assumption of normal error is favoured
 430 compared to lognormal error, so proceed with the results obtained
 431 from the NLR fit.
- 432 • If $\text{AIC}_{\text{cNLR}} - \text{AIC}_{\text{cLLR}} > 2$, the assumption of lognormal error is favoured
 433 compared to normal error, so proceed with the results obtained from
 434 the LLR fit.
- 435 • If $|\text{AIC}_{\text{cNLR}} - \text{AIC}_{\text{cLLR}}| \leq 2$, no model is favoured, so proceed with model

436 averaging:

$$437 \quad B_{\text{av}} = w_{\text{NLR}} V_{\text{NLR}} + w_{\text{LLR}} V_{\text{LLR}}$$

$$438 \quad \beta_{\text{av}} = w_{\text{NLR}} \beta_{\text{NLR}} + w_{\text{LLR}} \beta_{\text{LLR}}$$

439 where:

$$440 \quad w_{\text{NLR}} = \frac{1}{1 + e^{\frac{1}{2}(\text{AIC}_{\text{cNLR}} - \text{AIC}_{\text{cLLR}})}}$$

$$441 \quad w_{\text{LLR}} = \frac{1}{1 + e^{\frac{1}{2}(\text{AIC}_{\text{cLLR}} - \text{AIC}_{\text{cNLR}})}}$$

442 which are obtained to fulfill the next condition: $w_{\text{NLR}} + w_{\text{LLR}} = 1$. The

443 CIs for B_{av} and β_{av} are to be generated by ordinary bootstrapping^{II}.

- 444 3. Assess the validity of the underlying statistical assumptions with diagnos-
 445 tic plots because while it is rare for all the assumptions to be fully satisfied
 446 by real-life data sets, major violations indicate the lack of appropriateness
 447 of the model and, thus, the potential invalidity of the results.

^{II}*cmplxcrunner* has available the next bootstrapping alternatives (28): ordinary, “Resampling Residuals” method, “Wild” method, and “Monte-Carlo” method.

448 Calculating the coefficient of determination

449 We think the best approach in this situation is to apply the generalized R^2 that,
450 for continuous models, was defined as (26):

$$451 \quad R^2 = 1 - \left(\frac{\mathcal{L}(0)}{\mathcal{L}(\hat{\theta})} \right)^{\frac{2}{n}}$$

452 where $\mathcal{L}(\hat{\theta})$ and $\mathcal{L}(0)$ denote the likelihoods of the fitted and the “null” model,
453 respectively, and n is the sample size. In terms of the loglikelihoods, the gen-
454 eralized coefficient of determination would be:

$$455 \quad R^2 = 1 - e^{-\frac{2}{n}(\log \mathcal{L}(\hat{\theta}) - \log \mathcal{L}(0))}$$

456 We have the likelihoods calculated from the previous section, but what about
457 the “null” models? We understand that they are the models with only the
458 intercept. So for the Gaussian additive error model:

$$459 \quad \mathcal{L}_{\text{norm}}(0) = \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma_{\text{NLRO}}^2}} \exp\left(-\frac{(y_i - \bar{y})^2}{2\sigma_{\text{NLRO}}^2}\right) \right]$$

460 So:

$$\begin{aligned} 461 \quad \log \mathcal{L}_{\text{norm}}(0) &= -\frac{n}{2} \log |2\pi\sigma_{\text{NLRO}}^2| - \frac{1}{2\sigma_{\text{NLRO}}^2} \sum_{i=1}^n (y_i - \bar{y})^2 \\ 462 \quad &= -\frac{n}{2} (\log |2\pi\sigma_{\text{NLRO}}^2| + 1) \end{aligned}$$

463 since $\sigma_{\text{NLR0}}^2 = \frac{1}{n} \sum (y_i - \bar{y})^2 = \frac{1}{n} \text{TSS}_{\text{NLR}}$. Now, coming back to the coefficient of
 464 determination, we have:

$$\begin{aligned} 465 \quad R_{\text{NLR}}^2 &= 1 - e^{\frac{2}{n}(\log \mathcal{L}_{\text{NLR}}(0) - \log \mathcal{L}_{\text{NLR}}(\hat{\theta}))} = 1 - \exp\left(\frac{\log(\text{RSS}_{\text{NLR}})}{\log(\text{TSS}_{\text{NLR}})}\right) = \\ 466 \quad &= 1 - \frac{\text{RSS}_{\text{NLR}}}{\text{TSS}_{\text{NLR}}} = 1 - \frac{\sum_{i=1}^n (y_i - V_{\text{NLR}} x_i^{\beta_{\text{NLR}}})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \end{aligned}$$

467 recovering the traditional expression for R^2 . Using the same approach for cal-
 468 culating R_{LLR}^2 , then:

$$469 \quad \mathcal{L}_{\log n}(0) = \prod_{i=1}^n \left[\frac{1}{y_i \sqrt{2\pi\sigma_{\text{LLR0}}^2}} \exp\left(-\frac{(\log |y_i| - \overline{\log |y|})^2}{2\sigma_{\text{LLR0}}^2}\right) \right]$$

470 So:

$$\begin{aligned} 471 \quad \log \mathcal{L}_{\log n}(0) &= -\frac{n}{2} \log |2\pi\sigma_{\text{LLR0}}^2| - \frac{1}{2\sigma_{\text{LLR0}}^2} \sum_{i=1}^n (\log |y_i| - \overline{\log |y|})^2 - \sum_{i=1}^n \log |y_i| \\ 472 \quad &= -\frac{n}{2} (\log |2\pi\sigma_{\text{LLR0}}^2| + 1) - \sum_{i=1}^n \log |y_i| \end{aligned}$$

473 since $\sigma_{\text{LLR0}}^2 = \frac{1}{n} \sum (\log |y_i| - \overline{\log |y|})^2 = \frac{1}{n} \text{TSS}_{\log n}$. Again, recalling the expres-
 474 sion for the generalized coefficient of determination, we have:

$$\begin{aligned} 475 \quad R_{\text{LLR}}^2 &= 1 - e^{\frac{2}{n}(\log \mathcal{L}_{\text{LLR}}(0) - \log \mathcal{L}_{\text{LLR}}(\hat{\theta}))} = 1 - \exp\left(\frac{\log(\text{RSS}_{\text{LLR}})}{\log(\text{TSS}_{\text{LLR}})}\right) = \\ 476 \quad &= 1 - \frac{\text{RSS}_{\text{LLR}}}{\text{TSS}_{\text{LLR}}} = 1 - \frac{\sum_{i=1}^n (\log |y_i| - \log |V_{\text{LLR}}| - \beta_{\text{LLR}} \log |x_i|)^2}{\sum_{i=1}^n (\log |y_i| - \overline{\log |y|})^2} \end{aligned}$$

477 **X-weighted power-law fit**

478 When fitting the power-law of std vs. mean, we can take into account that
 479 every mean has uncertainty and estimate it for a sample size n by the SEM
 480 (*Standard Error of the Mean*):

$$481 \quad \text{SEM} = \frac{s}{\sqrt{n}}$$

482 where s is the sample standard deviation. So, the vector of weights is computed
 483 with:

$$484 \quad \mathbf{w} = \frac{1}{\text{SEM}} = \frac{\sqrt{n}}{s}$$

485 Here, the uncertainties affect the independent variable, so the fit is not so triv-
 486 ial as a Y-weighted fit, where the uncertainties affect the dependent variable.
 487 A standard approach to do this fit is: a) invert your variables before apply-
 488 ing the weights, b) then perform the weighted fit, and finally, c) revert the
 489 inversion. This method is deterministic, but the approximate solution worsens
 490 with smaller R^2 . For comparison, we develop a stochastic method by using a
 491 bootstrapping-like strategy that avoids the inversion and is applicable regard-
 492 less of R^2 . Both methods, detailed below, are implemented in *cmplxcruncher*.

493 **Method 1: By inverting the data**

494 In the case of the log-LR model, we have:

$$495 \quad \log y = \log V + \beta \log x \quad \rightarrow \quad \underbrace{\log x}_{\tilde{y}} = \underbrace{-\frac{1}{\beta} \log V}_b + \underbrace{\frac{1}{\beta}}_m \underbrace{\log y}_{\tilde{x}}$$

496 where m determines the slope or gradient of the fitted line, and b determines
497 the point at which the line crosses the y-axis, otherwise known as the y-intercept.

498 Once the model is fitted, the original parameters can be retrieved easily:

$$499 \quad \beta = \frac{1}{m}$$

$$500 \quad V = e^{-\beta b} = e^{-\frac{b}{m}}$$

501 Their respective uncertainties are to be obtained using *error propagation*:

$$502 \quad \sigma_\beta = \left| \frac{d\beta}{dm} \right| \sigma_m = \frac{1}{m^2} \sigma_m$$

$$503 \quad \sigma_V = \sqrt{\left(\frac{\partial V}{\partial b} \right)^2 \sigma_b^2 + \left(\frac{\partial V}{\partial m} \right)^2 \sigma_m^2} = \frac{1}{m} e^{-\frac{b}{m}} \sqrt{\sigma_b^2 + \frac{b^2}{m^2} \sigma_m^2}$$

504 **Method 2: Bootstrapping-like strategy**

505 The basic idea of bootstrapping is that inference about a population from sam-
506 ple data (sample \rightarrow population) can be modeled by resampling the sample
507 data and performing inference on (resample \rightarrow sample). To adapt this general

idea to our problem, we resample the x-data array using its errors array. That is, for each replicate, a new x-data array is computed based on:

$$x_i^* = x_i + v_i$$

where v_i is a Gaussian random variable with mean $\mu_i = 0$ and standard deviation $\sigma_i = \text{SEM}_i$, as defined previously. For each replicate a complete unweighted power-law fit is performed, as described in the previous section. It is worth mentioning that each replicate is filtered to avoid values of x_i^* under `eps` (obtained by `np.finfo(np.double).eps`) in order to keep away from the error of getting log of negatives or zero during the fit.

We devised and implemented a multi-step algorithm to estimate the fit parameters that finishes when a relative error of less than 10^{-4} is achieved. It also ends if the number of steps reaches 100 to avoid too much time lapse, to prevent any pathologic numeric case which, in fact, we still have not detected in all the data sets analyzed.

In the previous version of the algorithm, for each step, the method generated 10 replicates for each x-data point, in other words, it was computing the fit for 10 times the length of the x-data array replicates, with a maximum of 10000 fits per step. Nevertheless, we found that such an approach depending on the length of the x-data array did not perform better, so we decided to simplify the method and fix the number of fits per step in 100. This latter approach improved the performance.

529 The parameters of the X-weighted fit are then estimated by averaging through
 530 all the replicate fits performed, and their errors are estimated by computing
 531 the standard deviation also for all the fits. At the end of each step, the relative
 532 error is calculated by comparing the fit parameters estimation in the last step
 533 with the previous one.

534 Finally, both the coefficient of determination of the fit and the coefficient of
 535 correlation between the fit parameters are estimated by averaging.

536 Rank Stability Index (RSI)

537 The Rank Stability Index is shown as a percentage in a separate bar on the
 538 right of the rank matrix plot provided by *cmplxcruncher*. The RSI is strictly
 539 1 for an element whose range never changes over time, and is strictly 0 for
 540 an element whose rank oscillates between the extremes from time to time. So,
 541 RSI is calculated, per element, as 1 less the quotient of the number of true rank
 542 hops taken between the number of maximum possible rank hops, all powered
 543 to p :

$$544 \quad \text{RSI} = \left(1 - \frac{\text{true rank hops}}{\text{possible rank hops}}\right)^p = \left(1 - \frac{D}{(N-1)(t-1)}\right)^p$$

545 where D is the total of rank hops taken by the studied element, N is the num-
 546 ber of elements that have been ranked, and t is the number of time samples.
 547 The power index p is arbitrarily chosen to increase the resolution in the stable
 548 region; the value in the current version of the code is $p = 4$.

Case	Condition	Colour
1	$1 \geq \text{RSI} > 0.99$	blue
2	$\text{RSI} > 0.90$	green
3	$\text{RSI} > 0.75$	orange
4	$\text{RSI} > 0.25$	red
5	$0.25 \geq \text{RSI} \geq 0$	black

Table 7. Colour code of the RSI percentage text shown in rank plots, following the first condition satisfied.

549 As an example of this “zooming” effect in the stable region, to match a linear
 550 ($p = 1$) RSI of 0.9 to a powered one of 0.1, we should select $p = 21.8543$. An
 551 alternative way to obtain this effect and exactly map a linear RSI of 0.9 to a
 552 non-linear RSI (RSI') of 0.1, is by applying the following function:

$$553 \quad \text{RSI}' = \frac{10^{10\left(1 - \frac{D}{(N-1)(t-1)}\right)} - 1}{10^{10} - 1} \approx 10^{-10\left(\frac{D}{(N-1)(t-1)}\right)}$$

554 where the approximation is valid because $10^{10} \gg 1$ but, the small price to pay
 555 for it is that, in the worst instability case, the RSI' would not be strictly 0 but
 556 10^{-10} .

557 The colour code of the RSI percentage text in the rank plot of *cmplxcruncher* is
 558 chosen following the first condition satisfied from those shown in Table 7 (see
 559 page 47).

560 Standardization

561 In order to show all the studies properly under common axes, we decided to
 562 standardize the Taylor parameters using the group of healthy individuals for
 563 each study. With this approach, all the studies can be visualized in a shared
 564 plot with units of Taylor-parameters standard-deviation on their axes.

565 For a Taylor parameter, e.g. V , the estimate of the mean (\hat{V}) for the healthy
 566 subpopulation, composed of h individuals, is:

$$567 \quad \hat{V} = \frac{1}{W_1} \sum_{i=1}^h V_i \omega_i = \sum_{i=1}^h V_i \omega_i$$

568 as $W_1 = \sum_i^h \omega_i = 1$, since ω_i are normalized weights calculated as:

$$569 \quad \omega_i = \frac{\frac{1}{\sigma_{V_i}^2}}{\sum_i^h \frac{1}{\sigma_{V_i}^2}}$$

570 being σ_{V_i} the estimation of the uncertainty in V_i obtained together with V_i from
 571 the X-weighted power-law fit described in Section , for healthy individuals.

572 Likewise, the estimation of the standard deviation for the healthy population
 573 ($\hat{\sigma}_V$) is:

$$574 \quad \hat{\sigma}_V = \sqrt{\frac{1}{W_1 - \frac{W_2}{W_1}} \sum_{i=1}^h [\omega_i (V_i - \hat{V})^2]}$$

575 being $W_2 = \sum_i^h \omega_i^2$, which finally yields to:

$$576 \quad \hat{\sigma}_V = \sqrt{\frac{1}{1 - \sum_i^h \omega_i^2} \sum_{i=1}^h [\omega_i (V_i - \hat{V})^2]}$$

577 Acknowledgments

578 Authors declare that there are no competing financial interests in relation to
579 the work described here.

580 Funding Information

581 This work is supported by Generalitat Valenciana Prometeo Grants II/2014/050,
582 II/2014/065, by the Spanish Grants FPA2011-29678, BFU2012-39816-C02-
583 01 of MINECO and by PITN-GA-2011-289442-INVISIBLES. JMM & DMM ac-
584 knowledge FPI and FISABIO fellowships. **Modificar becas de JMM y DMM**
585 **¿poner algún grant más?**

References

1. Qin J, Li Y, Cai Z, Li S, Zhu J, Zhang F, Liang S, Zhang W, Guan Y, Shen D, Peng Y, Zhang D, Jie Z, Wu W, Qin Y, Xue W, Li J, Han L, Lu D, Wu P, Dai Y, Sun X, Li Z, Tang A, Zhong S, Li X, Chen W, Xu R, Wang M, Feng Q, Gong M, Yu J, Zhang Y, Zhang M, Hansen T, Sanchez G, Raes J, Falony G, Okuda S, Almeida M, LeChatelier E, Renault P, Pons N, Batto J-M, Zhang Z, Chen H, Yang R, Zheng W, Li S, Yang H, Wang J, Ehrlich SD, Nielsen R, Pedersen O, Kristiansen K, Wang J. 2012. A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature* **490**:55–60.
2. Brown JM, Hazen SL. 2015. The Gut Microbial Endocrine Organ: Bacterially Derived Signals Driving Cardiometabolic Diseases. *Annu Rev Med* **66**:343–359.
3. Durbán A, Abellán JJ, Jiménez-Hernández N, Artacho A, Garrigues V, Ortiz V, Ponce J, Latorre A, Moya A. 2013. Instability of the faecal microbiota in diarrhoea-predominant irritable bowel syndrome. *FEMS Microbiol Ecol* **86**:581–589.
4. Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, Schwager E, Knights D, Song SJ, Yassour M, Morgan XC, Kostic AD, Luo C, González A, McDonald D, Haberman Y, Walters T, Baker S, Rosh J, Stephens M, Heyman M, Markowitz J, Baldas-

- sano R, Griffiths A, Sylvester F, Mack D, Kim S, Crandall W, Hyams J, Huttenhower C, Knight R, Xavier RJ. 2014. The treatment-naive microbiome in new-onset Crohn's disease. *Cell Host Microbe* **15**:382–392.
5. Ridaura VK, Faith JJ, Rey FE, Cheng J, Duncan AE, Kau L, Griffi NW, Lombard V, Henrissat B, Bain JR, Michael J, Ilkayeva O, Semenkovich CF, Funai K, Hayashi DK, Lyle J, Martini MC, Ursell LK, Clemente JC, Treuren W Van, William A, Knight R, Newgard CB, Heath AC, Gordon JI, Kau AL, Griffin NW, Muehlbauer MJ. 2013. Gut Microbiota from Twins Discordant for Obesity Modulate Metabolism in Mice Gut Microbiota from Twins Metabolism in Mice. *Science* **341**:1241214.
6. Turnbaugh PJ, Hamady M, Yatsunenko T, Cantarel BL, Duncan A, Ley RE, Sogin ML, Jones WJ, Roe BA, Affourtit JP, Egholm M, Henrissat B, Heath AC, Knight R, Gordon JI. 2009. LETTERS A core gut microbiome in obese and lean twins. *Nature* **457**:480–484.
7. Subramanian S, Huq S, Yatsunenko T, Haque R, Mahfuz M, Alam MA, Benezra A, DeStefano J, Meier ME, Muegge BD, Barratt MJ, VanArendonk LG, Zhang Q, Province MA, Petri WA, Ahmed T, Gordon JI. 2014. Persistent gut microbiota immaturity in malnourished Bangladeshi children. *Nature* **510**:417–21.
8. Caporaso, J.G. et al. Moving pictures of the human microbiome. *Genome Biol.* **12**, R50 (2011).

- 629 9. Faith, J.J. et al. The long-term stability of the human gut microbiota.
630 *Science* **341**, 1237439 (2013).
- 631 10. Smith M.I. et al. Gut microbiomes of Malawian twin pairs discordant for
632 kwashiorkor. *Science* **339**, 548-54 (2013).
- 633 11. David, L.A. et al. Diet rapidly and reproducibly alters the human gut
634 microbiome. *Nature* **505**, 559-63 (2014).
- 635 12. Dethlefsen L., Relman D. A. Incomplete recovery and individualized re-
636 sponses of the human distal gut microbiota to repeated antibiotic per-
637 turbation. *Proc. Nat. Acad. Sci. USA* **108**, 4554-61 (2011).
- 638 13. Caporaso, J.G. et al. QIIME allows analysis of high-throughput commu-
639 nity sequencing data. *Nature Methods* **7**, 335-6 (2010).
- 640 14. Ames, S.K. et al. Scalable metagenomic taxonomy classification using a
641 reference genome database. *Bioinformatics* **29**, 2253-2260 (2013).
- 642 15. Eislér, Z., Bartos, I., Kertész, J. Fluctuation scaling in complex systems:
643 Taylor's law and beyond. *Adv. Phys.* **57**, 85 (2008).
- 644 16. Taylor, L.R. Aggregation, Variance and the mean. *Nature* **189**, 732-35
645 (1961).
- 646 17. Jorgensen, B., Martinez, J.R., Tsao, M. Asymptotic behaviour of the vari-
647 ance function. *Scand. J. Statist.* **21**, 223-243 (1994).

- 648 18. Fronczak,A., Fronczak,P. Origins of Taylor’s power law for fluctuation
649 scaling in complex systems. *Phys. Rev. E* **81**, 066112 (2010).
- 650 19. Kendal, W.S., Jorgensen,B. Taylor’s power law and fluctuation scaling
651 explained by a central-limit-like convergence. *Phys. Rev. E* **83**, 066115
652 (2011).
- 653 20. Kendal, W.S., Jorgensen,B. Tweedie convergence: A mathematical basis
654 for Taylor’s power law. *Phys. Rev. E* **84**, 066120 (2011).
- 655 21. Blumm, N. et al. Dynamics of ranking processes in complex systems.
656 *Phys. Rev. Lett.* **109**, 128701 (2012).
- 657 22. Weber, J. *et al.* Fluctuation dissipation theorem. *Phys. Rev.* **101**, 1620-6
658 (1956).
- 659 23. Gordon, A., Hannon, G.J. FASTX-Toolkit. FASTQ/A shortreads pre-
660 processing tools (2010). http://hannonlab.cshl.edu/fastx_toolkit/ (ac-
661 cessed 23 Feb 2015).
- 662 24. Quast C. *et al.* The SILVA ribosomal RNA gene database project: im-
663 proved data processing and web-based tools (2013)
- 664 25. Xiao Xiao, Ethan P. White, Mevin B. Hooten, and Susan L. Durham. On
665 the use of log-transformation vs. nonlinear regression for analyzing bi-
666 ological power laws. *Ecology* **92**, 10, 1887-1894 (2011).

667 26. Magee L., R^2 measures based on wald and likelihood ratio joint signifi-
668 cance tests. *The American Statistician* **44**, 3, 250-253 (1990).

669 27. Nagelkerke N.J.D., A note on a general definition of the coefficient of
670 determination. *Biometrika* **78**, 3, 691-692 (1991).

671 28. Wu, C.F.J. Jackknife, bootstrap and other resampling methods in regres-
672 sion analysis. (with discussions) *The Annals of Statistics* **14**: 1261-1350
673 (1986)

674 Eliminar et al. y poner la referencia completa como exige la guía
675 de estilo de la revista...