

Distances between histograms

Renato Fabbri^{1, a)}*São Carlos Institute of Physics, University of São Paulo (IFSC/USP), PO Box 369, 13560-970, São Carlos, SP, Brazil*

(Dated: 28 October 2015)

This document presents reference values for a distance metric derived from the Kolmogorov-Smirnov statistical test. Each measure is a distance between two histograms. The sections are self-explanatory on deriving benchmarks by comparing samples from usual distributions and on exemplifying the power of the acquired knowledge.

PACS numbers: 05.10-a,

Keywords: Kolmogorov-Smirnov test, benchmark, distance measure, histogram

I. INTRODUCTION

Be $F_{1,n}$ and $F_{2,n'}$ two empirical cumulative distributions, where n and n' are the number of observations on each sample. The two-sample Kolmogorov-Smirnov test rejects the null hypothesis (that the histograms are the outcome of the same underlying distribution) if:

$$D_{n,n'} > c(\alpha) \sqrt{\frac{n+n'}{nn'}} \quad (1)$$

where $D_{n,n'} = \sup_x [F_{1,n} - F_{2,n'}]$ and $c(\alpha)$ are related to the critical region α by:

α	0.1	0.05	0.025	0.01	0.005	0.001
$c(\alpha)$	1.22	1.36	1.48	1.63	1.73	1.95

If distributions are drawn from empirical data, $D_{n,n'}$ is given as are n and n' . All terms in equation 1 are positive and $c(\alpha)$ can be isolated:

$$c(\alpha) < \frac{D_{n,n'}}{\sqrt{\frac{n+n'}{nn'}}} = c' \quad (2)$$

When c' is high, low values of α favor rejecting the null hypothesis. For example, when c' is greater than ≈ 1.7 , one might assume that $F_{1,n}$ and $F_{2,n'}$ are outcomes of different distributions. More importantly for us is that c' is a measure of distance between both distributions¹. The main contribution of the following sections is the explicit display of reference values from which one might derive knowledge from collections of empirical measures of c' or even of a single value of c' .

A. Philosophical and technological note

Difference and equivalence is of central role in human cognition, philosophy and science. This fact is so deeply

recognized that thinkers often reduce thought to classifications, e.g. through the mathematical concept of equivalence classes^{2,3}. Histograms are very immediate and informative wherever there is a phenomenon of interest which can yield measurements. This present document should enable conclusions to be drawn about the equivalence (and difference) of the processes underlying sets of measurements for a very broad range of phenomena. The minimum delivered by the following tables is that the software implementation that render them present measures of how different are the underlying distributions of given samples.

B. Document outline

Section II expose reference values drawn from simulations. Section III exemplifies the use of such reference values to make sense of phenomena. Section IV hold final remarks. Software and data specification are given in Appendix A.

II. REFERENCES THROUGH SIMULATIONS

On every case, values of c' are given for simulations involving at least normal, uniform, triangular, Weibull and power-law distributions.

A. When the null hypothesis is true

If the null hypothesis is true, than the number of rejections of the null hypothesis (that is: $c' > c(\alpha)$) in N_c comparisons should not exceed αN_c . To verify this, let $C' = \{c'_i\}$ be a set of c' measures, and $C'(\alpha) = \{c' : c' > c(\alpha)\}$. Be $|C'(\alpha)|$ the cardinality of $C'(\alpha)$, i.e. the number of comparisons in which the two-sample Kolmogorov-Smirnov test rejects the null hypothesis for a given α .

The most important result of this section is that $|C'(\alpha)|$ rarely exceeds αN_c , where N_c is the number of

^{a)} <http://ifsc.usp.br/~fabbri/>; Electronic mail: fabbri@usp.br

comparisons made, no matter what probability distribution type or the specific setting. Also important are that $c' > c(\alpha)$ in many cases and that α is a good estimate of the upper limit of the probability of such an event.

αN_c	α	$c(\alpha)$	$ C'_1(\alpha) $	$ C'_2(\alpha) $	$ C'_3(\alpha) $
10.0	0.100	1.22	2	3	5
5.0	0.050	1.36	1	1	3
2.5	0.025	1.48	0	1	1
1.0	0.010	1.63	0	1	0
0.5	0.005	1.73	0	0	0
0.1	0.001	1.95	0	0	0

TABLE I. The theoretical maximum number αN_c of rejections of the null hypothesis for critical values of α . The number of comparisons is $N_c = 100$, each with the sample size of $N_o = 1000$ observations. Each histogram have $N_b = 30$ equally spaced bins. The N_o values of c'_1 were calculated using simulations of normal distributions with $\mu = 0$ and $\sigma = 1$. The N_o values of c'_2 were calculated using simulations of normal distributions with $\mu = 3$ and $\sigma = 2$. The N_o values of c'_3 were calculated using simulations of normal distributions with $\mu = 6$ and $\sigma = 3$. Over all N_c comparisons, $\mu(c'_1) = 0.7015$ and $\sigma(c'_1) = 0.2315$, $\mu(c'_2) = 0.7176$ and $\sigma(c'_2) = 0.2658$, $\mu(c'_3) = 0.7348$ and $\sigma(c'_3) = 0.2527$.

αN_c	α	$c(\alpha)$	$ C'_1(\alpha) $	$ C'_2(\alpha) $	$ C'_3(\alpha) $
10.0	0.100	1.22	7	6	6
5.0	0.050	1.36	4	1	2
2.5	0.025	1.48	1	1	2
1.0	0.010	1.63	1	1	1
0.5	0.005	1.73	1	0	1
0.1	0.001	1.95	0	0	1

TABLE II. The theoretical maximum number αN_c of rejections of the null hypothesis for critical values of α . The number of comparisons is $N_c = 100$, each with the sample size of $N_o = 1000$ observations. Each histogram have $N_b = 30$ equally spaced bins. The N_o values of c'_1 were calculated using simulations of uniform distributions within $[0, 1)$. The N_o values of c'_2 were calculated using simulations of uniform distributions within $[2, 6)$. The N_o values of c'_3 were calculated using simulations of normal distributions with $\mu = 4$ and $\sigma = 10$. Over all N_c comparisons, $\mu(c'_1) = 0.7415$ and $\sigma(c'_1) = 0.2904$, $\mu(c'_2) = 0.7954$ and $\sigma(c'_2) = 0.2462$, $\mu(c'_3) = 0.7934$ and $\sigma(c'_3) = 0.2981$.

B. When the null hypothesis is false

The $m(c')$ are median values of c' . $\overline{C'(\alpha)} = \frac{|C'(\alpha)|}{N_c}$ is the fraction of rejection of the null hypothesis with critical region α .

αN_c	α	$c(\alpha)$	$ C'_1(\alpha) $	$ C'_2(\alpha) $	$ C'_3(\alpha) $	$ C'_4(\alpha) $
10.0	0.100	1.22	0	4	2	7
5.0	0.050	1.36	0	2	1	3
2.5	0.025	1.48	0	1	1	1
1.0	0.010	1.63	0	0	0	0
0.5	0.005	1.73	0	0	0	0
0.1	0.001	1.95	0	0	0	0

TABLE III. The theoretical maximum number αN_c of rejections of the null hypothesis for critical values of α . The number of comparisons is $N_c = 100$, each with the sample size of $N_o = 1000$ observations. Each histogram have $N_b = 30$ equally spaced bins. The N_o values of c'_1 were calculated using simulations of 1-parameter Weibull distributions with $a = 0.1$. The N_o values of c'_2 were calculated using simulations of 1-parameter Weibull distributions with $a = 2$. The N_o values of c'_3 were calculated using simulations of 1-parameter Weibull distributions with $a = 4$. Over all N_c comparisons, The N_o values of c'_4 were calculated using simulations of 1-parameter Weibull distributions with $a = 6$. Over all N_c comparisons, $\mu(c'_1) = 0.0635$ and $\sigma(c'_1) = 0.0388$, $\mu(c'_2) = 0.7178$ and $\sigma(c'_2) = 0.2499$, $\mu(c'_3) = 0.7167$ and $\sigma(c'_3) = 0.2379$, $\mu(c'_4) = 0.7073$ and $\sigma(c'_4) = 0.2583$.

III. EXAMPLE USES IN EMPIRICAL DATA

IV. CONCLUSIONS

ACKNOWLEDGMENTS

Financial support was obtained from CNPq (140860/2013-4, project 870336/1997-5), United Nations Development Program (contract: 2013/000566; project BRA/12/018) and FAPESP. The authors are grateful to the American Jewish Committee for maintaining an online copy of the Adorno book used on the epigraph⁷, to GMANE creators and maintainers for the public email list data, to the communities of the email lists and other groups used in the analysis, and to the Brazilian Presidency of the Republic for keeping Participabr code and data open. We are also grateful to developers and users of Python scientific tools.

Appendix A: Software and data specifications

⁷R. Chicheportiche and J.-P. Bouchaud, “Weighted kolmogorov-smirnov test: Accounting for the tails,” Physical Review E **86**, 041115 (2012).

σ	$\mu(c')$	$\sigma(c')$	$m(c')$	$\min(c')$	$\max(c')$	$C'(0.1)$	$C'(0.05)$	$C'(0.025)$	$C'(0.01)$	$C'(0.005)$	$C'(0.001)$
0.5	3.867	0.313	3.868	3.198,3.287,3.332	4.584,4.651,4.852	1.000	1.000	1.000	1.000	1.000	1.000
0.6	2.917	0.229	2.918	2.370,2.460,2.504	3.376,3.376,3.734	1.000	1.000	1.000	1.000	1.000	1.000
0.7	2.232	0.284	2.191	1.476,1.744,1.789	2.817,2.996,3.153	1.000	1.000	0.990	0.990	0.990	0.840
0.8	1.579	0.279	1.576	0.917,1.051,1.051	2.191,2.214,2.370	0.910	0.790	0.600	0.450	0.250	0.110
0.9	1.000	0.241	0.962	0.581,0.581,0.671	1.610,1.677,1.722	0.170	0.090	0.050	0.020	0.000	0.000
1.0	0.691	0.253	0.648	0.268,0.291,0.335	1.319,1.364,1.364	0.050	0.020	0.000	0.000	0.000	0.000
1.1	0.936	0.274	0.906	0.470,0.492,0.492	1.565,1.565,1.610	0.170	0.080	0.070	0.000	0.000	0.000
1.2	1.304	0.252	1.297	0.783,0.827,0.894	1.878,1.990,2.080	0.620	0.370	0.210	0.100	0.070	0.020
1.3	1.773	0.276	1.766	1.029,1.319,1.342	2.326,2.348,2.504	0.990	0.970	0.840	0.680	0.540	0.250
1.4	2.126	0.259	2.080	1.565,1.677,1.677	2.683,2.750,2.795	1.000	1.000	1.000	0.990	0.960	0.740
1.5	2.422	0.285	2.404	1.856,1.901,1.923	2.996,3.086,3.130	1.000	1.000	1.000	1.000	1.000	0.940
1.6	2.777	0.287	2.806	2.012,2.191,2.191	3.399,3.399,3.466	1.000	1.000	1.000	1.000	1.000	1.000
1.7	3.070	0.264	3.086	2.504,2.527,2.571	3.667,3.667,3.690	1.000	1.000	1.000	1.000	1.000	1.000
1.8	3.363	0.273	3.365	2.728,2.728,2.795	3.935,4.047,4.249	1.000	1.000	1.000	1.000	1.000	1.000
1.9	3.635	0.277	3.645	3.019,3.063,3.130	4.137,4.181,4.584	1.000	1.000	1.000	1.000	1.000	1.000
2.0	3.849	0.263	3.801	3.309,3.332,3.421	4.472,4.539,4.562	1.000	1.000	1.000	1.000	1.000	1.000

TABLE IV. Location and dispersion of $N_c = 100$ measurements of c' through simulations with normal distributions and $N_o = 1000$ events each. $N_b = 30$ equal bins were used to make the histograms. One normal distribution is fixed, with $\mu = 0$ and $\sigma = 1$, and compared against normal distributions with $\mu = 0$ and different values of σ .

μ	$\mu(c')$	$\sigma(c')$	$m(c')$	$\min(c')$	$\max(c')$	$C'(0.1)$	$C'(0.05)$	$C'(0.025)$	$C'(0.01)$	$C'(0.005)$	$C'(0.001)$
0.0	0.718	0.245	0.693	0.313,0.335,0.358	1.342,1.364,1.386	0.050	0.020	0.000	0.000	0.000	0.000
0.1	1.206	0.454	1.140	0.470,0.492,0.514	2.102,2.258,2.415	0.400	0.350	0.300	0.250	0.160	0.060
0.2	2.028	0.429	2.001	1.207,1.230,1.252	2.862,2.929,3.063	0.990	0.940	0.900	0.810	0.720	0.530
0.3	2.774	0.399	2.750	1.856,2.012,2.012	3.690,3.757,3.824	1.000	1.000	1.000	1.000	1.000	0.990
0.4	3.719	0.439	3.712	2.817,2.885,2.996	4.584,4.606,4.651	1.000	1.000	1.000	1.000	1.000	1.000
0.5	4.539	0.433	4.573	3.488,3.600,3.667	5.344,5.389,5.434	1.000	1.000	1.000	1.000	1.000	1.000
0.6	5.374	0.479	5.367	4.070,4.137,4.450	6.350,6.418,6.529	1.000	1.000	1.000	1.000	1.000	1.000
0.7	6.258	0.454	6.261	5.098,5.121,5.277	7.111,7.133,7.200	1.000	1.000	1.000	1.000	1.000	1.000
0.8	6.992	0.450	6.909	5.903,6.172,6.261	7.893,7.983,8.005	1.000	1.000	1.000	1.000	1.000	1.000
0.9	7.838	0.453	7.804	6.865,6.887,6.909	8.765,8.944,9.034	1.000	1.000	1.000	1.000	1.000	1.000
1.0	8.616	0.433	8.665	7.625,7.692,7.804	9.347,9.481,9.749	1.000	1.000	1.000	1.000	1.000	1.000

TABLE V. Location and dispersion of $N_c = 100$ measurements of c' through simulations with normal distributions and $N_o = 1000$ events each. $N_b = 30$ equal bins were used to make the histograms. One normal distribution is fixed, with $\mu = 0$ and $\sigma = 1$, and compared against normal distributions with different values of μ and fixed $\sigma = 1$.

b	$\mu(c')$	$\sigma(c')$	$m(c')$	$\min(c')$	$\max(c')$	$C'(0.1)$	$C'(0.05)$	$C'(0.025)$	$C'(0.01)$	$C'(0.005)$	$C'(0.001)$
0.7	6.748	0.360	6.731	6.015,6.060,6.104	7.714,7.714,7.737	1.000	1.000	1.000	1.000	1.000	1.000
0.75	5.465	0.340	5.456	4.740,4.852,4.875	6.149,6.373,6.440	1.000	1.000	1.000	1.000	1.000	1.000
0.8	4.484	0.288	4.494	3.712,3.779,3.846	5.031,5.031,5.031	1.000	1.000	1.000	1.000	1.000	1.000
0.85	3.304	0.264	3.287	2.795,2.862,2.885	3.913,3.913,3.980	1.000	1.000	1.000	1.000	1.000	1.000
0.9	2.313	0.273	2.337	1.722,1.766,1.789	2.795,2.929,3.153	1.000	1.000	1.000	1.000	0.990	0.900
0.95	1.331	0.261	1.364	0.760,0.850,0.872	1.878,2.012,2.124	0.630	0.510	0.280	0.110	0.040	0.020
1.0	0.767	0.270	0.716	0.358,0.358,0.402	1.386,1.476,1.588	0.080	0.050	0.010	0.000	0.000	0.000
1.05	1.305	0.275	1.297	0.805,0.850,0.850	1.878,1.901,2.147	0.570	0.410	0.240	0.120	0.080	0.010
1.1	2.158	0.244	2.158	1.565,1.677,1.677	2.683,2.773,2.885	1.000	1.000	1.000	0.990	0.970	0.810
1.15	3.026	0.263	2.996	2.437,2.527,2.549	3.533,3.555,3.846	1.000	1.000	1.000	1.000	1.000	1.000
1.2	3.686	0.270	3.690	3.019,3.130,3.175	4.271,4.271,4.360	1.000	1.000	1.000	1.000	1.000	1.000
1.25	4.499	0.311	4.472	3.891,3.891,3.913	5.232,5.389,5.568	1.000	1.000	1.000	1.000	1.000	1.000
1.3	5.113	0.317	5.121	4.360,4.405,4.427	5.724,5.724,5.881	1.000	1.000	1.000	1.000	1.000	1.000
1.35	5.772	0.304	5.780	5.165,5.188,5.210	6.350,6.418,6.596	1.000	1.000	1.000	1.000	1.000	1.000

TABLE VI. Location and dispersion of $N_c = 100$ measurements of c' through simulations with uniform distributions and $N_o = 1000$ events each. $N_b = 30$ equal bins were used to make the histograms. One uniform distribution has the fixed domain $[0, 1)$. The other uniform distribution in each comparison is also centered around 0.5, but spread over $b = b_u - b_l$ there b_l and b_u are the lower and upper boudaries.

a	$\mu(c')$	$\sigma(c')$	$m(c')$	$\min(c')$	$\max(c')$	$C'(0.1)$	$C'(0.05)$	$C'(0.025)$	$C'(0.01)$	$C'(0.005)$	$C'(0.001)$
0.01	0.029	0.012	0.022	0.022,0.022,0.022	0.067,0.067,0.089	0.000	0.000	0.000	0.000	0.000	0.000
0.1	0.139	0.100	0.112	0.022,0.022,0.022	0.402,0.402,0.492	0.000	0.000	0.000	0.000	0.000	0.000
0.3	1.660	0.627	1.621	0.402,0.470,0.514	2.840,2.862,3.511	0.730	0.660	0.580	0.500	0.420	0.350
0.5	3.988	0.427	3.991	2.549,2.683,2.974	4.763,4.785,4.942	1.000	1.000	1.000	1.000	1.000	1.000
0.7	3.448	0.522	3.365	2.482,2.639,2.728	4.696,4.964,5.523	1.000	1.000	1.000	1.000	1.000	1.000
0.9	3.016	0.347	2.974	2.326,2.393,2.415	3.734,3.779,3.846	1.000	1.000	1.000	1.000	1.000	1.000
1.1	1.923	0.309	1.889	1.364,1.431,1.453	2.639,2.706,2.907	1.000	1.000	0.950	0.820	0.680	0.420
1.3	1.252	0.274	1.207	0.604,0.738,0.805	1.834,1.945,2.057	0.480	0.340	0.190	0.090	0.050	0.010
1.5	0.693	0.240	0.682	0.291,0.291,0.335	1.207,1.275,1.498	0.020	0.010	0.010	0.000	0.000	0.000
1.7	1.078	0.289	1.062	0.581,0.626,0.648	1.856,1.923,2.012	0.240	0.130	0.080	0.070	0.040	0.010
1.9	1.586	0.337	1.576	0.559,0.939,0.962	2.303,2.393,2.460	0.880	0.780	0.590	0.440	0.320	0.140
2.1	2.130	0.331	2.102	1.431,1.588,1.588	2.795,2.817,2.907	1.000	1.000	0.990	0.960	0.850	0.670
2.3	2.587	0.331	2.605	1.923,1.990,2.012	3.287,3.354,3.421	1.000	1.000	1.000	1.000	1.000	0.990
2.5	3.155	0.358	3.153	2.214,2.281,2.303	3.757,3.935,3.958	1.000	1.000	1.000	1.000	1.000	1.000
2.7	3.557	0.368	3.600	2.817,2.885,2.885	4.517,4.517,4.651	1.000	1.000	1.000	1.000	1.000	1.000
2.9	3.884	0.425	3.868	2.907,2.952,3.108	4.718,4.763,4.942	1.000	1.000	1.000	1.000	1.000	1.000

TABLE VII. Location and dispersion of $N_c = 100$ measurements of c' through simulations with 1-parameter Weibull distributions and $N_o = 1000$ events each. $N_b = 30$ equal bins were used to make the histograms. One Weibull distribution has the fixed shape parameter $a = 1.5$. The other Weibull distribution in each comparison has varied values of a .