

# Воспроизведение результатов статьи в [py\\_graphs](#).

Владимир Ивашкин

5 июля 2018 г.

## 2 Logarithmic vs. plain measures

Здесь RI, а не ARI

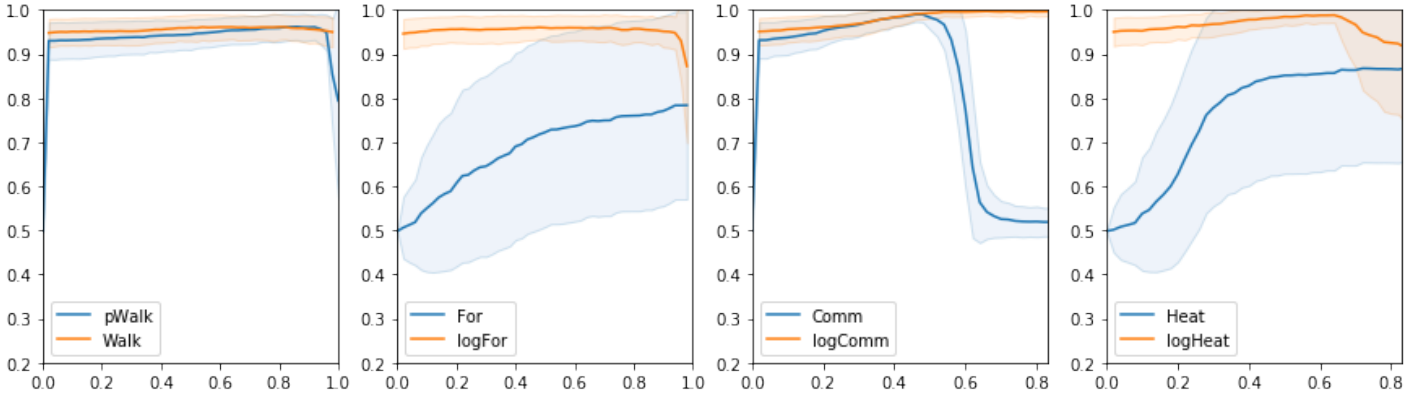


Рис. 1: Logarithmic vs. plain measures for  $G(100, (2)0.2, 0.05)$

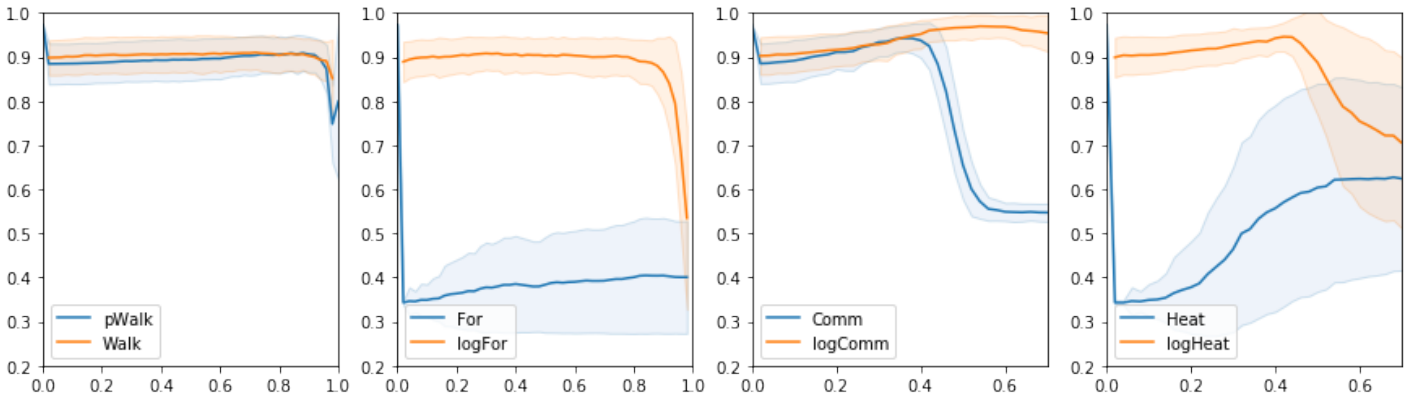


Рис. 2: Logarithmic vs. plain measures for  $G(100, (3)0.3, 0.1)$

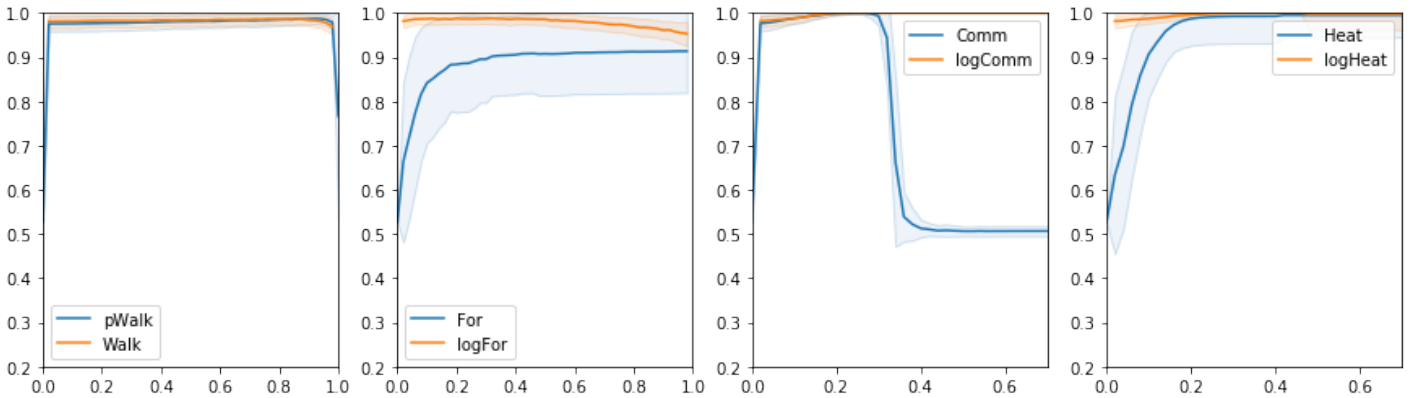


Рис. 3: Logarithmic vs. plain measures for  $G(200, (2)0.3, 0.1)$

### 3 Competition by Copeland's score

<b>Nodes</b>	100	100	100	100	200	200	200	200	<b>Sum</b>
<b>Classes</b>	2	2	4	4	2	2	4	4	<b>of</b>
<b><math>p_{out}</math></b>	0.1	0.15	0.1	0.15	0.1	0.15	0.1	0.15	<b>scores</b>
logComm	383	547	476	-66	301	565	592		<b>2798</b>
Comm	249	150	308	418	291	212	325		<b>1953</b>
SCCT	316	299	166	44	293	392	412		<b>1922</b>
logHeat	308	314	180	-264	301	321	343		<b>1503</b>
pWalk	-81	26	56	418	-105	-155	6		<b>165</b>
SCT	-74	36	78	44	47	-24	44		<b>151</b>
Heat	221	-342	-456	418	295	205	-478		<b>-137</b>
RSP	-96	4	62	-272	-32	-85	-30		<b>-449</b>
Walk	-90	-26	78	-222	-149	-125	-26		<b>-560</b>
logFor	-92	-44	-24	-264	-63	-92	-32		<b>-611</b>
FE	-202	-64	-44	-224	-135	-169	-134		<b>-972</b>
For	-387	-566	-456	418	-525	-574	-478		<b>-2568</b>
SP-CT	-455	-334	-424	-448	-519	-471	-544		<b>-3195</b>

(a) optimal parameters

<b>Nodes</b>	100	100	100	100	200	200	200	200	<b>Sum</b>
<b>Classes</b>	2	2	4	4	2	2	4	4	<b>of</b>
<b><math>p_{out}</math></b>	0.1	0.15	0.1	0.15	0.1	0.15	0.1	0.15	<b>scores</b>
logComm	413	568	448	356	332	568	598	598	<b>3881</b>
SCCT	269	274	136	78	340	391	423	360	<b>2271</b>
logHeat	318	183	290	142	340	273	202	98	<b>1846</b>
Comm	168	151	222	172	286	258	333	178	<b>1768</b>
SCT	58	92	46	90	26	45	38	104	<b>499</b>
logFor	-114	60	56	110	-55	-115	4	88	<b>34</b>
Walk	-84	-10	132	86	-140	-85	30	66	<b>-5</b>
pWalk	-125	-40	54	74	-163	-79	-2	-14	<b>-295</b>
FE	-198	-27	-27	120	-120	-186	-66	32	<b>-472</b>
RSP	-151	-1	-8	78	-138	-179	-106	-16	<b>-521</b>
Heat	299	-341	-502	-490	340	154	-417	-515	<b>-1472</b>
SP-CT	-463	-345	-320	-228	-558	-462	-446	-396	<b>-3218</b>
For	-390	-564	-588	-588	-490	-583	-591	-583	<b>-4377</b>

(b) 90th percentiles

Таблица 1: Copeland's scores of the measure families on random graphs

## 4 Reject curves

Measure (kernel)	$G(100, (2)0.3, 0.05)$ Opt. parameter, ARI	$G(100, (2)0.3, 0.1)$ Opt. parameter, ARI	$G(100, (2)0.3, 0.15)$ Opt. parameter, ARI
pWalk	0.93, 1.00	0.87, 0.91	0.73, 0.66
Walk	0.93, 1.00	0.67, 0.91	0.70, 0.65
For	0.60, 0.99	0.97, 0.51	0.40, 0.01
logFor	0.70, 1.00	0.40, 0.93	0.10, 0.68
Comm	0.33, 1.00	0.33, 0.98	0.30, 0.77
logComm	0.33, 1.00	0.47, <b>1.00</b>	0.57, <b>0.91</b>
Heat	0.37, 1.00	0.60, 0.87	0.73, 0.15
logHeat	0.37, 1.00	0.53, 0.99	0.37, 0.80
SCT	0.40, 1.00	0.57, 0.94	0.43, 0.72
SCCT	0.03, 1.00	0.57, 0.98	0.63, 0.80
RSP	0.97, 1.00	0.97, 0.93	0.97, 0.67
FE	0.90, 1.00	0.90, 0.91	0.87, 0.68
SP-CT	0.00, 0.99	0.03, 0.78	0.07, 0.49

Таблица 2: Optimal family parameters and the corresponding ARI's

Ошибка была в том, что подобранные параметры из таблицы выше принадлежат к диапазону  $[0, 1]$ , а значит их нужно преобразовывать к диапазону, специфичному для конкретной метрики. Я же этого не делал. Вторая ошибка состояла в том, что я использовал тут близости вместо расстояний. Еще тогда, когда я строил их в прошлый раз, я заметил, что по близостям logComm совсем не обгоняет остальные меры, но по расстояниям эффект выраженный. Тут его тоже видно:

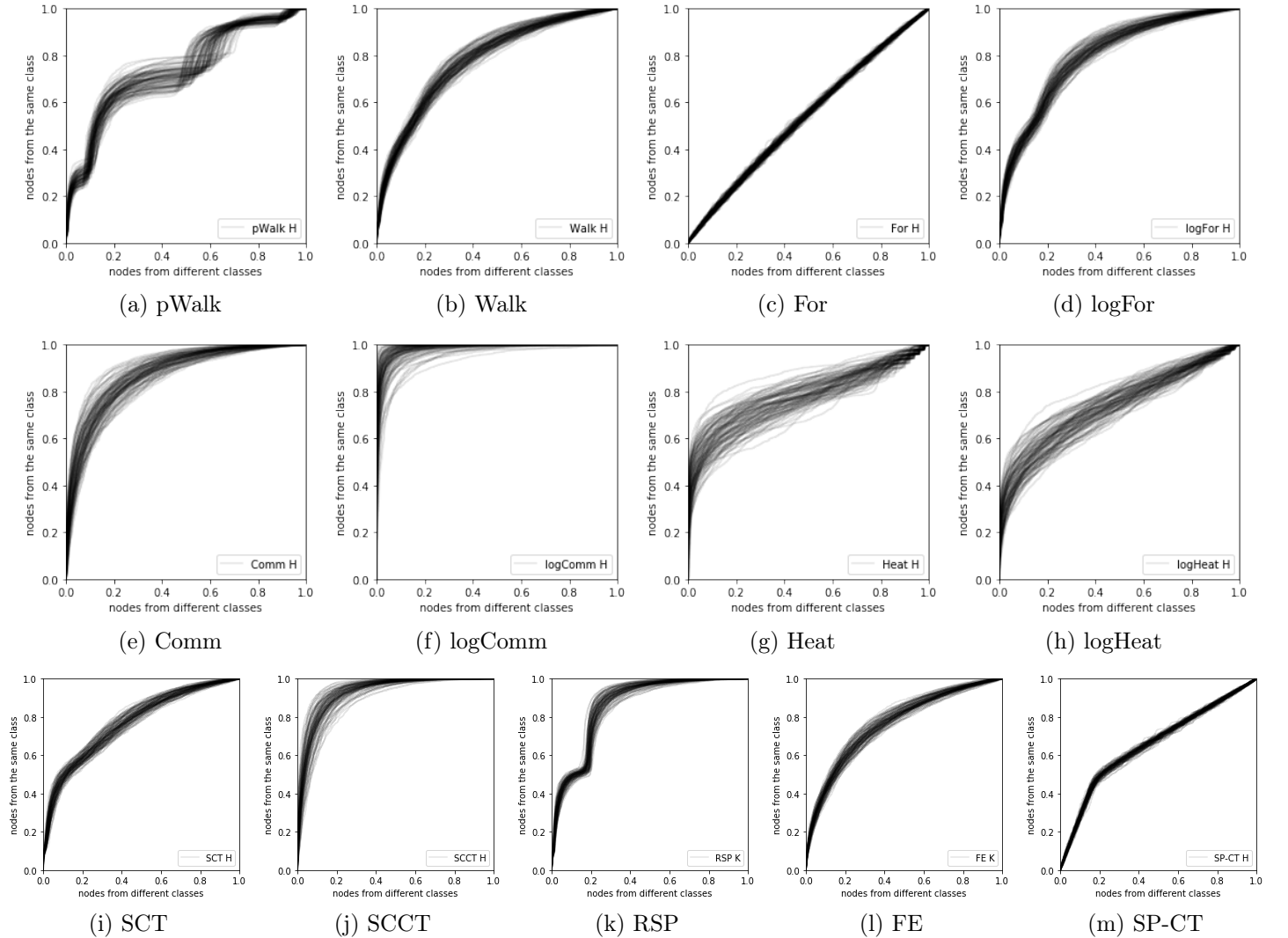
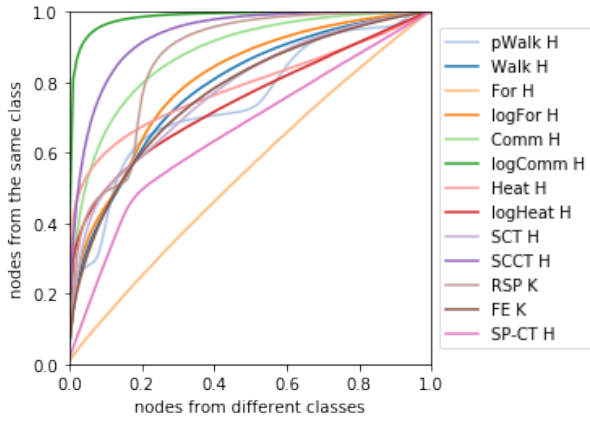
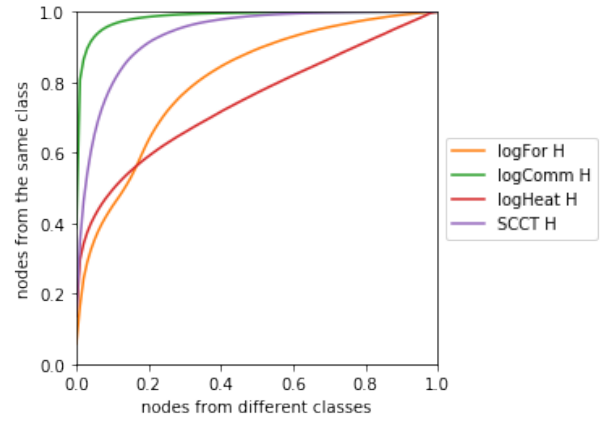


Рис. 4: Reject curves for the graph measures under study



(a) All families



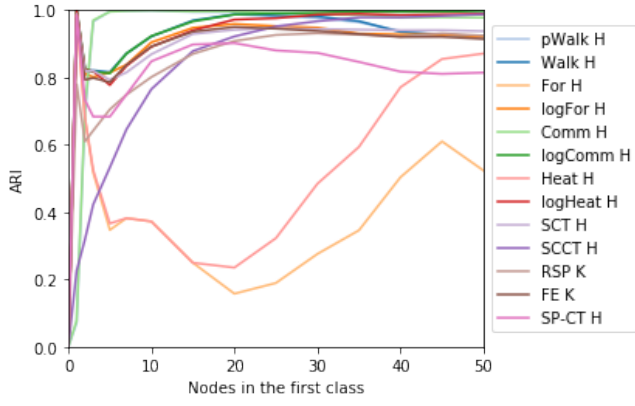
(a) All families

Рис. 5: Average reject curves

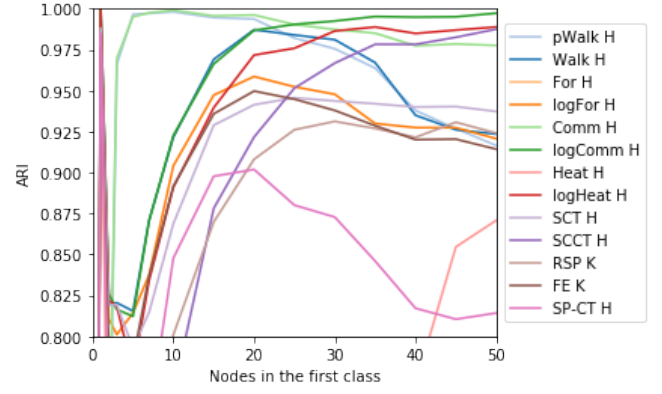
Здесь была проблема со взятием корня из Comm/logComm. А проблема была такая: если некоторые значения матрицы  $D$  при взятии корня превращаются в nan, то стандартная сортировка оставляет их на тех же позициях и отдельно сортирует массив слева и справа от них. Получается кусочно-возрастающая функция, из которой потом получаются несколько маленьких reject curve вместо одной большой. Решение – фильтровать эти nan и сортировать без них. Раз такой эффект вообще возник, значит в матрице  $D$  иногда появляются отрицательные значения.

Можно подозревать внешний вид графика pWalk. Может быть, это связано с тем, как мы фиксируем параметр. Параметром считаем отскалированное в  $[0, 1]$  число, для каждого графа преобразуем его в зависимости от спектрального радиуса матрицы  $A$  ( $param = t/\rho(A)$ ),  $t \in [0, 1]$ ).

## 5 Graphs with classes of different sizes



(a) All families



(b) Leading families

Рис. 6: Graphs with two classes of different sizes: clustering with optimal parameter values

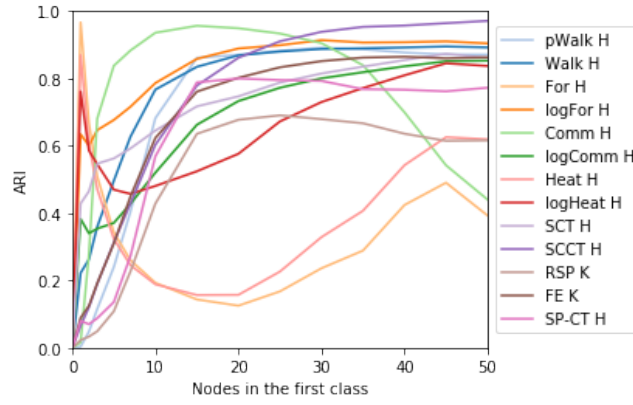
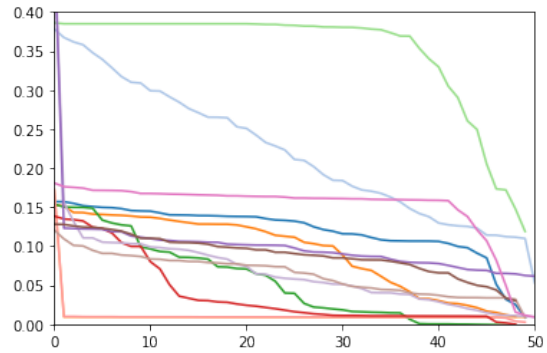


Рис. 7: Graphs with two classes of different sizes: random parameter values

$$P = \begin{pmatrix} 0.30 & 0.20 & 0.10 & 0.15 & 0.07 & 0.25 \\ 0.20 & 0.24 & 0.08 & 0.13 & 0.05 & 0.17 \\ 0.10 & 0.08 & 0.16 & 0.09 & 0.04 & 0.12 \\ 0.15 & 0.13 & 0.09 & 0.20 & 0.02 & 0.14 \\ 0.07 & 0.05 & 0.04 & 0.02 & 0.12 & 0.04 \\ 0.25 & 0.17 & 0.12 & 0.14 & 0.04 & 0.40 \end{pmatrix}.$$



of various measure families on a structure with 6 classes

.45.45ARI

## 6 Cluster analysis on several classical datasets

Здесь ошибка была в том, что я зафиксировал число классов – 2, хотя в датасете football их 12. Теперь все похоже на статью:

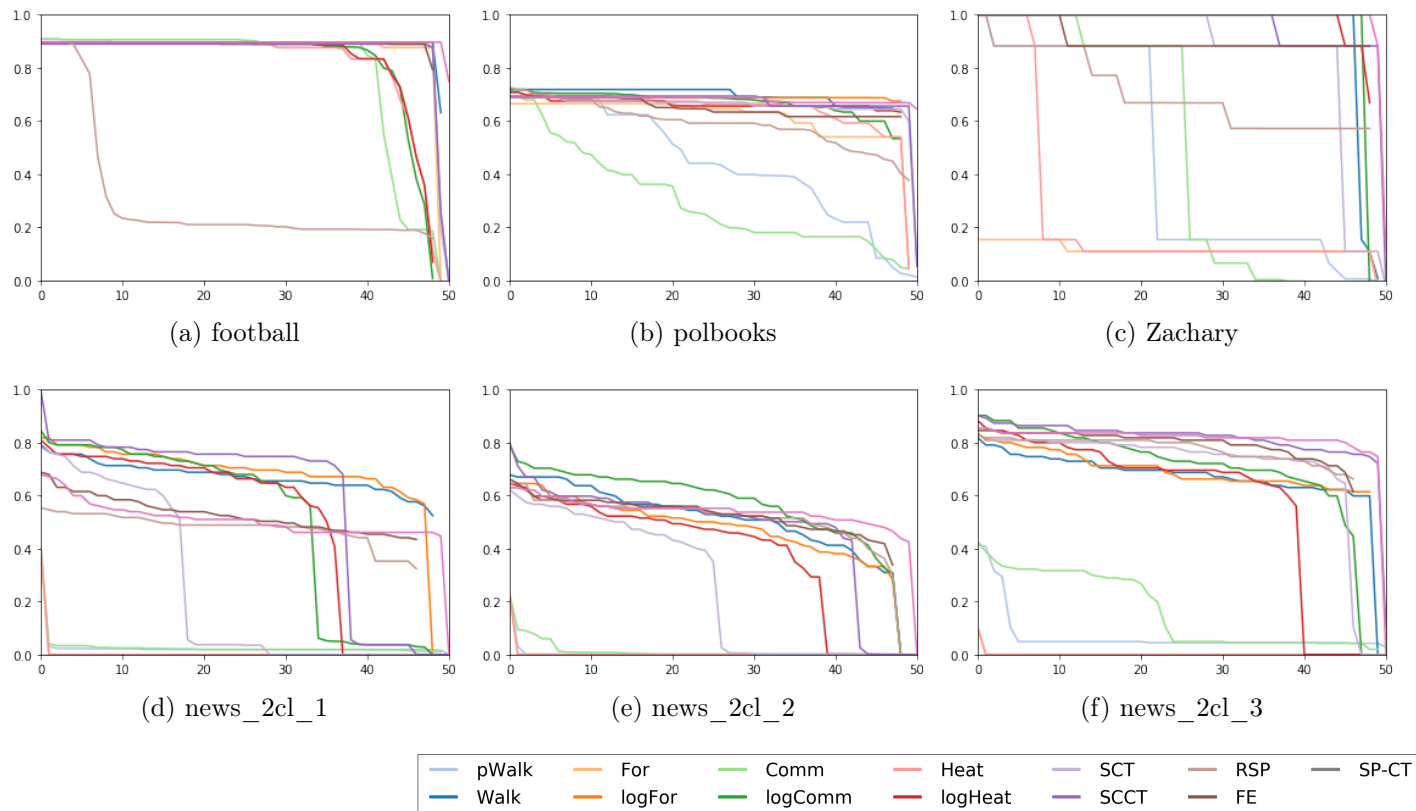


Рис. 8: ARI of various measure families on classical datasets