

# Finalised versions of normalisation and threshold

## Contents

<b>1</b>	<b>Idea</b>	<b>1</b>
<b>2</b>	<b>Computing distances</b>	<b>1</b>
<b>3</b>	<b>Normalize distances</b>	<b>2</b>
<b>4</b>	<b>Does normalisation work?</b>	<b>2</b>
<b>5</b>	<b>Choose thresholds for harmonies</b>	<b>3</b>
<b>6</b>	<b>Does it match with the threshold?</b>	<b>4</b>
<b>7</b>	<b>Should they tally with each other?</b>	<b>4</b>
7.1	cricket data . . . . .	4

## 1 Idea

Even after excluding clashes, the list of harmonies left could be large and overwhelming for human consumption. Hence, there is a need to rank the harmonies basis how well they capture the variation in the measured variable and additionally reduce the number of harmonies for further exploration/visualization. Gestalt theory suggests that when items are placed in close proximity, people assume that they are in the same group because they are close to one another and apart from other groups. Hence, displays that capture more variation within different categories in the same group would be important to bring out different patterns of the data. Thus the idea here is to rate a harmony pair higher if this variation between different levels of the x-axis variable is higher on an average across all levels of facet variables.

## 2 Computing distances

One of the potential ways to evaluate this variation is by computing the pairwise distances between the distributions of the measured variable. We do this through Jensen-Shannon distance which is based on Kullback-Leibler divergence. Probability distributions are represented through sample quantiles instead of kernel density estimate so that there is minimal dependency on selecting kernel or bandwidth.

We shall call this measure of variation as Median Maximum Pairwise Distances (MMPD)

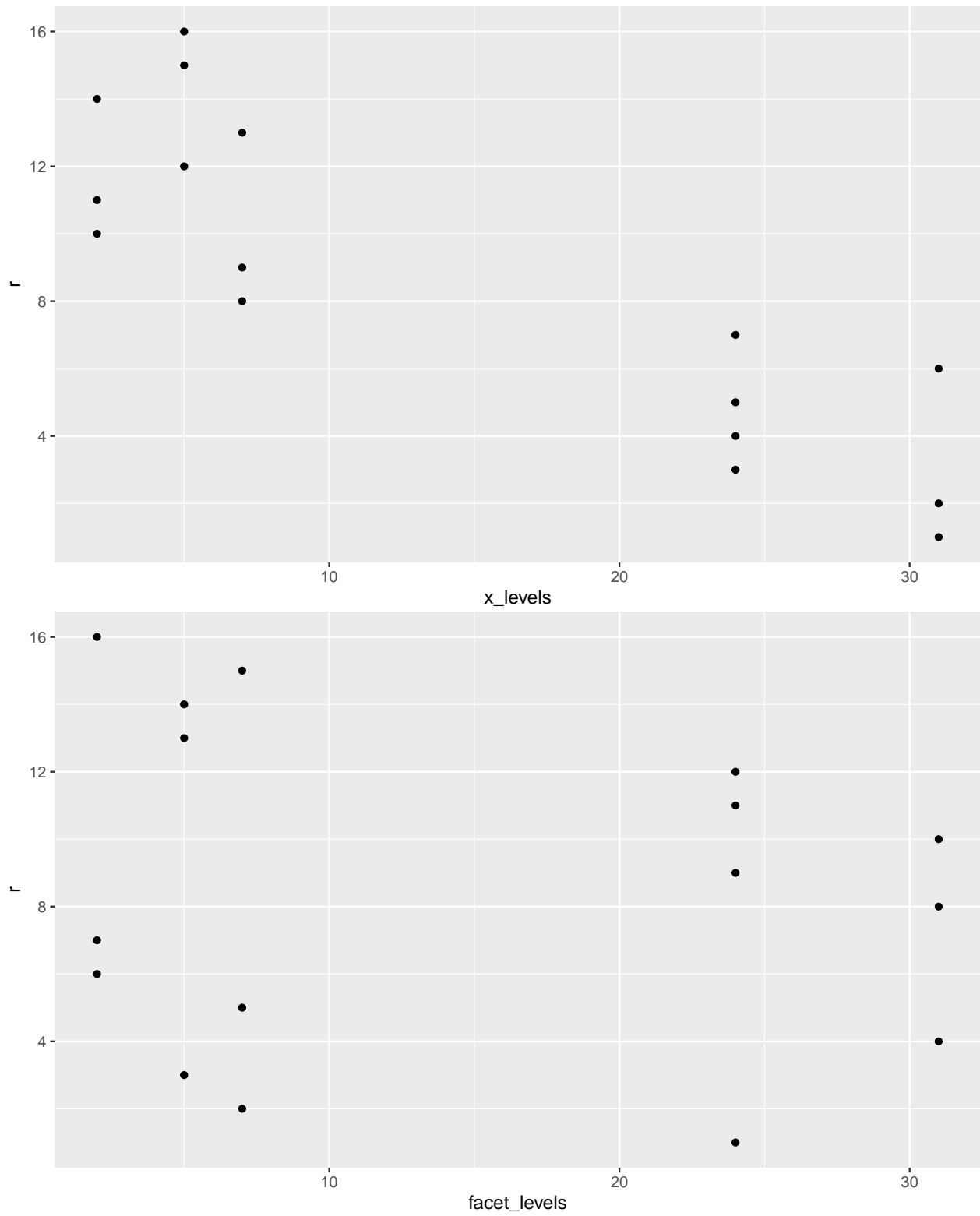
### 3 Normalize distances

The harmony pairs could be arranged from highest to lowest average maximum pairwise distances across different levels of the harmonies. But maximum is not robust to the number of levels and is higher for harmonies with higher levels. Thus these maximum pairwise distances need to be normalized for different harmonies in a way that eliminates the effect of different levels. The Fisher–Tippett–Gnedenko theorem in the field of Extreme Value Theory states that the maximum of a sample of iid random variables after proper re-normalization can converge in distribution to only one of Weibull, Gumbel or Freschet distribution, independent of the underlying data or process. The normalizing constants, however, vary depending on the underlying distribution and hence it is important to assume a distribution of distances in our case.

### 4 Does normalisation work?

(Show with by comparing with maximum and showing for similar levels)

facet_variable	x_variable	facet_levels	x_levels	MMPD	max_pd	r
wknd_wday	hour_day	2	24	0.72071	0.63930	7
wknd_wday	day_month	2	31	0.62645	0.70867	6
day_week	day_month	7	31	0.41627	2.21756	2
week_month	hour_day	5	24	0.24523	1.25759	3
day_week	hour_day	7	24	0.19640	0.99914	5
hour_day	day_month	24	31	0.17692	2.34599	1
day_month	hour_day	31	24	0.13035	1.12679	4
day_month	day_week	31	7	0.01898	0.35166	8
day_month	wknd_wday	31	2	0.01359	0.06925	10
hour_day	wknd_wday	24	2	0.00879	0.03487	11
week_month	wknd_wday	5	2	0.00684	0.02082	14
hour_day	day_week	24	7	0.00453	0.09317	9
week_month	day_week	5	7	0.00347	0.02765	13
day_week	week_month	7	5	0.00226	0.01999	15
wknd_wday	week_month	2	5	0.00219	0.00176	16
hour_day	week_month	24	5	0.00195	0.03016	12



## 5 Choose thresholds for harmonies

Permutation test:

**Assumption:** random permutation without considering ordering (global)

1. Given the data;  $\{v_t : t = 0, 1, 2, \dots, T-1\}$ , the MMPD is computed and is represented by  $MMPD_{obs}$ .
2. From the original sequence a random permutation is obtained:  $\{v_t^* : t = 0, 1, 2, \dots, T-1\}$ .
3. MMPD is computed for all random permutation of the data and is represented by  $MMPD_{sample}$ .
4. Steps (2) and (3) are repeated a large number of times M (e.g. 1000).
5. For each permutation, one  $MMPD_{sample}$  value is obtained.
6. 95<sup>th</sup> percentile of this  $MMPD_{sample}$  distribution is computed and stored in  $MMPD_{threshold}$ .
7. If  $MMPD_{obs} > MMPD_{threshold}$ , harmony pairs are accepted. Only one threshold for all harmony pairs.

Pros: Considering thresholds global for all harmony pairs would imply less computation time.

Cons: Only one threshold for all harmony pairs might not be an appropriate measure but a good benchmark.

## 6 Does it match with the threshold?

## 7 Should they tally with each other?

### 7.1 cricket data

