Step 4 Normalisation

19th JRC Annual training on Composite Indicators and Scoreboards

Marcos Domínguez-Torreiro



10 STEPS to build a Composite Indicator





Normalisation

• What is it?

It is the process of transforming the (raw) values of the indicators into a meaningful common numerical scale

- What are we doing it for?
 - To produce numerical scores that help us make direct comparisons of performance across indicators
 - As a prior step to compute composite indicators, i.e. as a prior step to build meaningful **aggregate/summary performance metrics** that combine the performance information from the indicators in the framework using **additive/multiplicative formulas** (weighted averages)
- When shall we apply a normalisation method to the indicators in the framework?

Whenever we have to adjust for the presence of:

- different units of measurement (e.g. euros vs. percentages)
- different distributions and/or variance (variablility) to prevent the undue influence of variables with larger variability on the aggregate scores – even when the unit of measurement is the same, numerical scales might not be comparable (e.g. GDP per capita expressed in euros vs. R+I funds per capita also expressed in euros)
- different "polarity" of the indicators' measurement scales (the higher the value the better vs the lower the value the better)

N.B.: Different normalisation methods may lead to different rankings – the impact of the alternative options can always be explored at the "Robustness & Sensitvity" analysis step

Main methods for normalising indicators

Outline

Min-max Distance to target/reference Z-scores

5. Ranking/percentile ranking

4. Categorical scale

6. Quantile normalization*

Non-linear transformations (reshape the original distributions)

Linear

transformations (do

not modify the shape of the original

distributions)



Example: employment rate (+) and long-term unemployment rate (-)



- The **long-term unemployment rate** is the share of the economically active population aged 15 to 74 who has been unemployed for 12 months or more.
- The **employment rate** is the percentage of employed persons in relation to the total population in the age group 20-64
- The indicators have **different polarity** (i.e. raw values **negatively correlated**).



Linear transformations

• A linear transformation formula creates a new variable by **adding (or subtracting) a constant and/or multiplying (dividing)** the variable **by another constant**.

e.g.1) $I = \frac{X - lgp}{ugp - lgp}$, where *lgp* is the "lower" reference goalpost and *ugp* the "upper" reference goalpost *e.g.2)* $I = \frac{X - avg}{sd}$, where *avg* is the mean and sd the standard deviation of the original variable

• Linear transformation formulas **do not modify the basic shape of the distribution** of the original variable:

- skewness and kurtosis values will not change
- it does not adjust for the presence of extreme values
- pairwise (Pearson) correlation coefficients between the indicators in the framework will not be affected

Min-max

$$I_j = \frac{X_j - lgp}{ugp - lgp} = \frac{X_j - \min(X)}{\max(X) - \min(X)}$$

For indicators with positive polarity (reverse polarity during normalisation using 1 - *I*)

Effects of normalisation:

- Unifies scale of measurement
- Unifies range: [0, 1]
- μ , σ^2 not equal
- Extreme values can be identified
- Shape of distribution: no adjustments



Min-max (pos.)



norm_emp

mean	var	min	max	range	skew	kurtosis
0.64	0.08	0.00	1.00	1.00	-0.92	0.02





kurtosis

0.02

80





Min-max (neg.)

	mean	var	min	max	range	skew	kurtosis
unemp_rate_LT	2.15	2.51	0.50	7.70	7.20	2.05	5.02









kurtosis

5.02



Min-max





r = 0.74



Relative distance to target/reference

 X_j : raw value country j

 X_R : raw reference value ("upper" goalpoast – internal or external reference point) m_R : minimum theoretical value

$$I_{j} = \frac{X_{j} - lgp}{ugp - lgp} = \frac{X_{j} - m_{R}}{X_{R} - m_{R}} \stackrel{\downarrow}{=} \frac{X_{j}}{X_{R}}$$

For indicators with positive polarity (reverse polarity during normalisation using $1 - I_i$)

Effects of normalisation:

- Unifies scale of measurement
- Range not unified
- μ , σ^2 not equal
- Extreme values can be identified
- Shape of distribution: no adjustments



Relative distance to target/reference (pos.)

	mean	var	min	max	range	skew	kurtosis
emp_rate	76.38	24.63	64.80	82.90	18.10	-0.92	0.02





 mean
 var
 min
 max
 range
 skew
 kurtosis

 norm_emp
 0.92
 0.00
 0.78
 1.00
 0.22
 -0.92
 0.02











Relative distance to target/reference (neg.)



Commission

Relative distance to target/reference







r = 0.74

Z-scores



For indicators with positive polarity (reverse polarity during normalisation multiplying Z_j by -1)

Effects of normalisation:

- Unifies scale of measurement
- Range not unified
- Unifies $\mu = 0$, $\sigma^2 = 1$
- Extreme values can be identified
- Shape of distribution: no adjustments



Z-scores (pos.)







Z-scores (neg.)

17

		mean	var	min n	nax rang	e skew	kurtosis
	unemp_rate_LT	2.15	2.51	0.50 7	.70 7.20) 2.05	5.02
° ° ° ° ° ° ° ° °		0 0	0.4 - 0.3 - 2002 - / 0.1 -				
AT BE BG CY CZ DE DK EE EL ES FI FR HR HU IE IT country_code	LT LU LV MT NL PL PT RO SE	si sk		2	4 unemp_rat	e_LT	ute e i e
	norm_unemp 0	.00 1.	.00 -3.8	50 1.04	4.54	-2.05 KU	5.02
° ° ° ° ° ° ° ° ° ° ° ° ° ° ° ° ° ° °	, ° ° °	0	0.6 -				\sim
• •		0	0.4 Atistep 0.2				
			0.0 -				





Z-scores







r = 0.74

Non-linear transformations

• Non-linear (ordinal) normalisation methods **transform the shape of the original distribution** of values.

• Non-linear transformations **impact the** (Pearson) **correlation metrics across indicators** in the normalised data set.

• Any **information regarding differences in levels** that might be present in the original indicators **is lost** in their normalised form.

• Normalised values are not affected by (not sensitive to) the presence of extreme values in the original data set, i.e. **extreme values can no longer be identified** (disappear) **in the transformed data set**.



Ranking

 $I_j = rank^d (X_j)$

Ranks are usually assigned in *descending order* due to ease of communication: for indicators with positive polarity the highest value would receive rank #1 and the lowest value rank #n (reverse polarity during normalisation using $n - rank^d (X_i)$)

Effects of normalisation:

- Unifies scale of measurement
- Unifies range: [1, n]
- Unifies μ , σ 2
- Extreme values cannot be identified
- Reshapes and unifies distributions

N.B: In the presence of missing values, percentile-ranks could be considered as an alternative to rankings: each raw value is replaced by the corresponding *p*-th percentile value to secure a unified range of normalised values (0, 100] in the transformed indicators.

N.B.: In normalised datasets where lower normalised values (ranks) represent better performance, ranks must be inverted (e.g. *n* – *rank*) before calculating *geometric averages* to avoid rewarding unbalanced performance across indicators. Once the geometric average is calculated, the transformation must be reversed (e.g. *n* – *geom.avg*) before interpreting the results.



Ranking (pos.)





Ranking (neg.)



Ranking





European Commission

r = 0.70

Categorical scales



(reverse polarity during normalisation multiplying *x* by -1)

Effects of normalisation:

- Unifies scale of measurement
- Unifies range
- Unifies μ , σ 2
- Extreme values cannot be identified
- Reshapes and unifies distributions



Categorical scales (pos.)





Categorical scales (neg.)

	mean	var	min	max	range	skew	kurtosis
unemp_rate_LT	2.15	2.51	0.50	7.70	7.20	2.05	5.02





	mean	var	min	max	range	skew	kurtosis
norm_unemp	51.85	900.28	0.00	100.00	100.00	-0.12	-0.74









Categorical scales



r = 0.69



Quantile normalisation*

• Assigns different scores to each observation in the indicator, while transforming original distributions into standard normal distributions.

• It applies a **two-stage normalisation method** that **first** converts the original distribution of values into a **uniform distribution**, and **second** transforms this uniform distribution into a **standard normal distribution** of scores, with mean zero and standard deviation equal to one.

• The two-stage normalisation formula is as follows:

$$I_j = \operatorname{qnorm}\left(\frac{\operatorname{rank}(X_j)}{n+1}\right)$$

where *qnorm()* is the function that returns the value matching a given cumulative probability under a normal distribution with mean zero and standard deviation one and *rank()* is a sorting function that assigns ranks in ascending order.



Key messages

	Lin	ear transformati	Non-linear transformations		
_		Distance to		Categorical	
	Min-max	reference	Ranks	scales	
Unifies scale of measurement	Y	Y	Y	Y	Y
Unifies range	Y	Ν	Ν	Y	Y
Unifies mean/variance	Ν	Ν	Y	Y	Y
Extreme values can be identified	Y	Y	Y	N	N
Reshapes and unifies distributions	Ν	Ν	Ν	Y	Y



Additional caveats

• When normalising, we must **pay attention to the nature of the original information** contained in the raw variables (see e.g. Stevens, 1946 *On the theory of scales of measurement*).

• If we are dealing with strictly qualitative ordinal-scale information, we should be careful about trying to make inferences about the extent of the actual gap in "real world" performance based on the magnitude of the differences in the numerals used in the raw/normalised scale of measurement.

• Raw/normalised variables fulfil the "**intervals**" **property** if we can safely assume that the numerals assigned to each observation are meaningful representations of actual differences in the real world, i.e. if equal differences in the magnitude of those numerals represent equal differences in the underlying characteristic or phenomenon.

• When a scale of measurement does not fulfil the intervals property, equal gaps in values (e.g. equal gaps in rank positions, or equal differences in Likert scale values) do not necessarily represent equal differences in performance as regards the underlying "real world" attribute or characteristic; therefore, under those circumstances descriptive statistics such as mean and variance become *numerical operations* instead of *meaningful statistics*.



Thank you



<u>marcos.dominguez-torreiro@ec.europa.eu | jrc-coin@ec.europa.eu</u> <u>https://knowledge4policy.ec.europa.eu/composite-indicators_en</u>

© European Union 2023

Unless otherwise noted the reuse of this presentation is authorised under the <u>CC BY 4.0</u> license. For any use or reproduction of elements that are not owned by the EU, permission may need to be sought directly from the respective right holders.

All images © European Union 2023

