# On the probability distribution of typological frequencies

## Michael Cysouw

Max Planck Institute for Evolutionary Anthropology, Leipzig

# VO vs. OV word order



Dryer, Matthew S. (2005) 'Order of object and verb' in Martin Haspelmath et al. (eds.)
*World Atlas of Language Structures* (Oxford: Oxford University Press), 338-41.

|  | Verb - Object | Object - Verb | Total |
| --- | --- | --- | --- |
| Noun - Relative Clause | **370** | **96** | 466 |
| Relative Clause - Noun | **5** | **109** | 114 |
| Total | 375 | 205 | 580 |

|  | Verb - Object | Object - Verb | Total |
|---|---|---|---|
| Noun - Relative Clause | **63.8%** | **16.6%** | 80.3% |
| Relative Clause - Noun | **0.9%** | **18.8%** | 19.7% |
| Total | 64.7% | 35.3% | 100% |

|  | Verb - Object | Object - Verb | Total |
|---|---|---|---|
| Noun - Relative Clause | + 3.96 | – 5.35 | |
| Relative Clause - Noun | – 8.00 | + 10.82 | |
| Total | | | |

"In a representative sample of languages, if no universal were involved, i.e. if the distribution of types along some parameter were purely random, then we would expect each type to have roughly an equal number of representatives. To the extent that the actual distribution departs from this random distribution, the linguist is obliged to state and, if possible, account for this discrepancy."

Comrie, Bernard (1989) *Language Universals and Linguistic Typology*, (2nd edition), Oxford: Blackwell).

What kind of random distributions are we dealing with?

# Size of Phoneme Inventory

- Lehfeldt (1975)
  **gamma distribution**

- Justeson & Stephens (1984)
  **log-normal distribution**

- Maddieson (2005)
  **normal distribution**

# Linguistic "Complexity"

- Nichols (1992)
  **normal distribution**

- Nichols et al. (2006)
  **"bell-shaped" distribution**

# All data from WALS



Shapiro-Wilk p = 0.02

Dryer, Matthew S. (2005) Order of subject, object and verb' in Martin Haspelmath et al. (eds.) *World Atlas of Language Structures* (Oxford: Oxford University Press), 330-33.

Genera

250

200

150

100

50

0

SOV    SVO    VSO    VOS    OVS    OSV

Dryer, Matthew S. (2005) Order of subject, object and verb' in Martin Haspelmath et al. (eds.) *World Atlas of Language Structures* (Oxford: Oxford University Press), 330-33.

# Noun Phrase word order

- "Those three red books"

- Demonstrative - Numeral - Adjective - Noun

- Dem-Num-A-N

Cinque, Guglielmo (2005) 'Deriving Greenberg's Universal 20 and
Its Exceptions' *Linguistic Inquiry* 36 (3), 315 - 332.
Dryer, Matthew S. (2006) 'On Cinque on Greenberg's universal 20'
Presentation at MPI-EVA Leipzig.

100

80

60

40

20

0

N-A-Num-Dem
Dem-Num-A-N
Num-N-A-Dem
N-A-Dem-Num
Dem-N-A-Num
Dem-Num-N-A
Dem-A-N-Num
N-Dem-A-Num
N-Num-A-Dem
N-Dem-Num-A
Num-N-Dem-A
A-N-Dem-Num
Num-A-N-Dem
Dem-N-Num-A
Num-Dem-N-A
N-Num-Dem-A
A-N-Num-Dem
Num-Dem-A-N
A-Dem-Num-N
A-Dem-N-Num
Dem-A-Num-N
Num-A-Dem-N
A-Num-Dem-N
A-Num-N-Dem

Why should we believe any claimed distribution ?

# Searching for distributions

- Empirical data won't help much

- Fitting data to a distribution always results in a more or less good fit.

- It is unclear what it would mean for an empirical distribution to have a "reasonably good fit"

We need some kind of theoretical notion about how typological distributions arise !

# Size of Phoneme Inventory
## according to Justeson & Stephens (1984)

- phonemes are based on features

- the number of phonemic features are normally distribution across languages (?!)

- with $n$ features one can make $2^n$ combinations

- thus, the sizes of phoneme inventories are log-normally distruted

# Variation in typological studies

- Which languages are investigated ?

- How many languages are investigated ?

- On what basis are the types distinguished ?

- How many different types are distinguished ?

Sample of
Linguist

Languages in
today's world

All human languages, past & present

# Markov-chain queueing model

- Languages change

- Being of a certain type is like being in a queue of a shop being set up by a linguist

- Language enter the queue following a poisson distribution, and they exit the queue according to a poisson distribution

- The length of the queue is then negatively exponentially distributed

# B. Wälchli's data on motion events

- 72 languages

- 335 clauses for each language from Bible

- lexical verbs describing motion events

|  | MRD | LIT | ENG | FRE |
|---|---|---|---|---|
| 1050 | sams | eiti | go | aller |
| 1070 | sams | eiti | come | venir |
| 1090 | sams | eiti | come | venir |
| 1104 | lisems | kopti | come | sortir |
| 1105 | valgoms | zengti | descend | descendre |
| 1114 | – | – | come | se faire entendre |
| 1120 | vetjams | varyti | drive | pousser |
| 1140 | sams | eiti | come | se rendre |
| 1160 | jutams | eiti | walk | marcher |

disagreement (d)

both different

agreement (a)

Jaccard similarity:
$$a/_{a+d}$$

|  | MRD | LIT | ENG | FRE |
|---|---|---|---|---|
| 1050 | sams | eiti | go | aller |
| 1070 | sams | eiti | come | venir |
| 1090 | sams | eiti | come | venir |
| 1104 | lisems | kopti | come | sortir |
| 1105 | valgoms | zengti | descend | descendre |
| 1114 | – | – | come | se faire entendre |
| 1120 | vetjams | varyti | drive | pousser |
| 1140 | sams | eiti | come | se rendre |
| 1160 | jutams | eiti | walk | marcher |

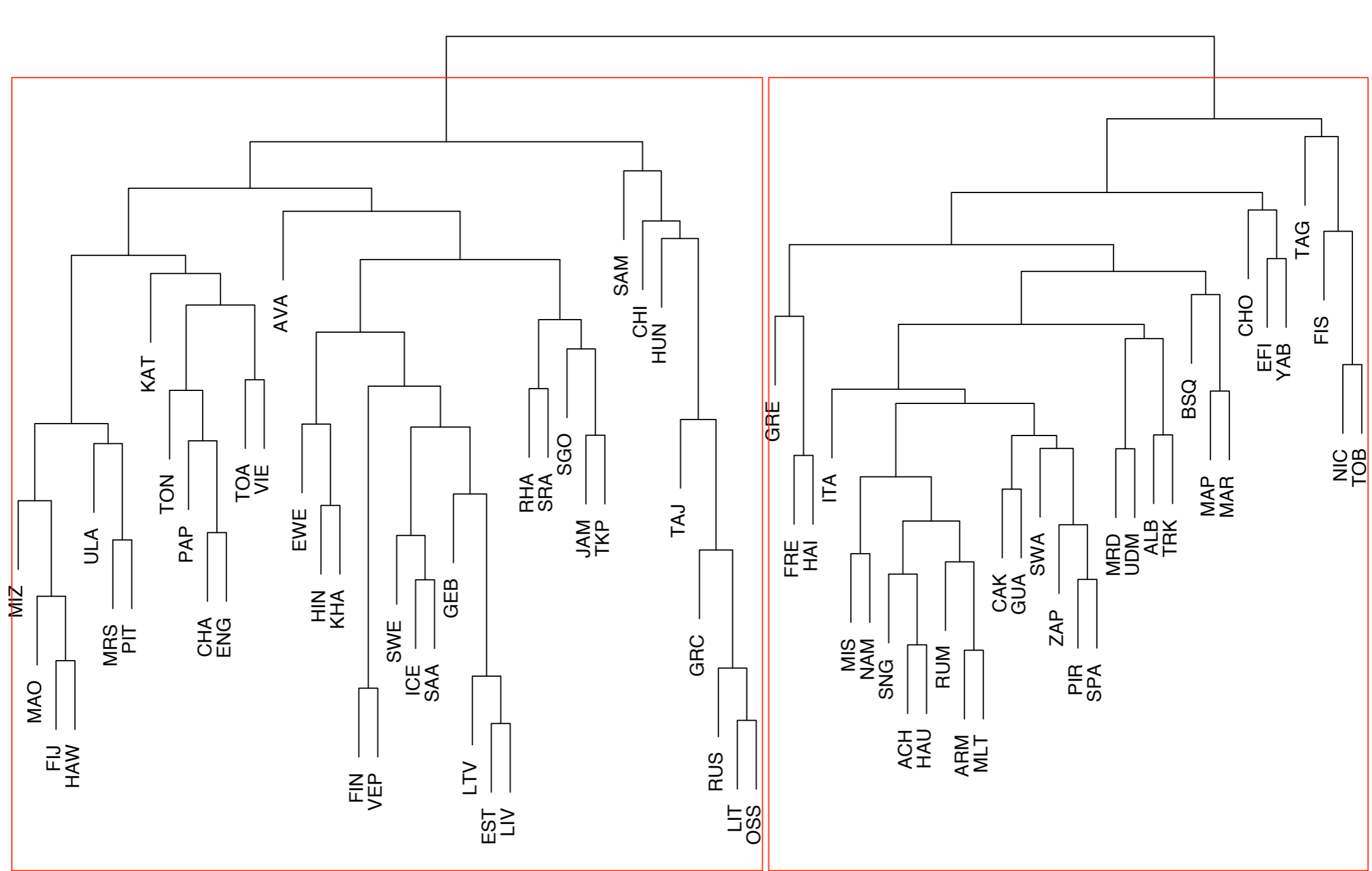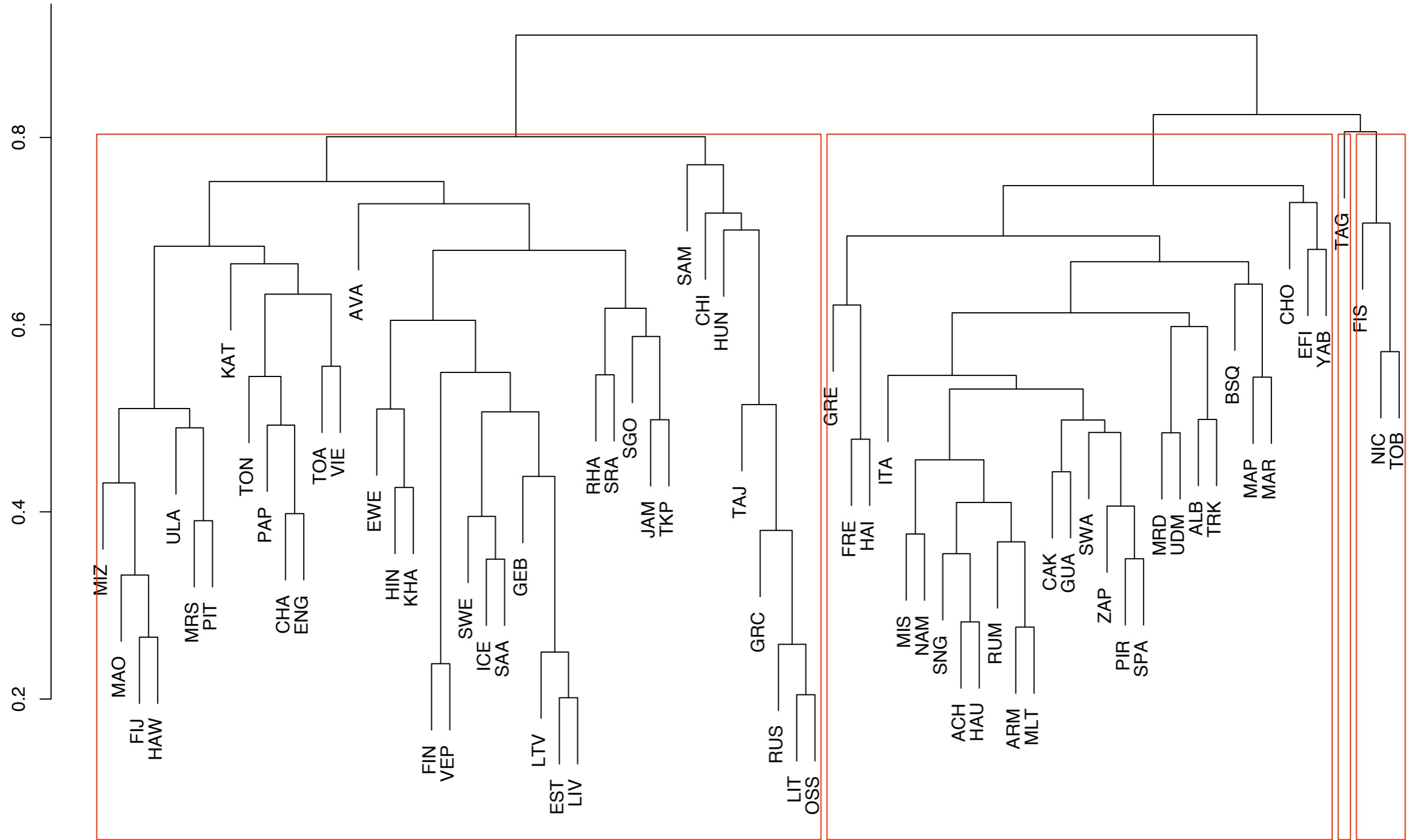| | ACH | ALB | ARM | AVA | BSQ | CAK | CHA | CHI | CHO | EFI | ENG | EST | EWE | FIJ | FIN | FIS | FRE | GEB | GRC | GRE | GUA | HAI | HAU | HAW | HIN | HUN | ICE | ITA | JAM | KAT | KHA | LIT | LIV | LTV | MAO | MAP | MA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ACH | 1 | 0.43 | 0.62 | 0.31 | 0.41 | 0.57 | 0.6 | 0.22 | 0.41 | 0.51 | 0.55 | 0.46 | 0.54 | 0.35 | 0.46 | 0.35 | 0.45 | 0.44 | 0.4 | 0.39 | 0.52 | 0.41 | 0.72 | 0.41 | 0.55 | 0.23 | 0.48 | 0.47 | 0.32 | 0.36 | 0.55 | 0.37 | 0.47 | 0.45 | 0.36 | 0.43 | 0.4 |
| ALB | 0.43 | 1 | 0.53 | 0.26 | 0.41 | 0.45 | 0.41 | 0.2 | 0.32 | 0.34 | 0.38 | 0.3 | 0.33 | 0.26 | 0.31 | 0.29 | 0.4 | 0.33 | 0.25 | 0.41 | 0.53 | 0.35 | 0.49 | 0.29 | 0.34 | 0.28 | 0.3 | 0.5 | 0.3 | 0.28 | 0.37 | 0.26 | 0.28 | 0.33 | 0.28 | 0.4 | 0.4 |
| ARM | 0.62 | 0.53 | 1 | 0.27 | 0.4 | 0.63 | 0.61 | 0.26 | 0.37 | 0.4 | 0.54 | 0.46 | 0.47 | 0.34 | 0.42 | 0.4 | 0.41 | 0.44 | 0.4 | 0.39 | 0.57 | 0.33 | 0.67 | 0.36 | 0.53 | 0.24 | 0.46 | 0.49 | 0.32 | 0.31 | 0.49 | 0.34 | 0.46 | 0.49 | 0.34 | 0.37 | 0.3 |
| AVA | 0.31 | 0.26 | 0.27 | 1 | 0.27 | 0.24 | 0.27 | 0.25 | 0.18 | 0.22 | 0.32 | 0.36 | 0.27 | 0.31 | 0.33 | 0.21 | 0.26 | 0.32 | 0.24 | 0.25 | 0.31 | 0.23 | 0.26 | 0.28 | 0.32 | 0.25 | 0.34 | 0.26 | 0.28 | 0.25 | 0.32 | 0.29 | 0.33 | 0.34 | 0.3 | 0.27 | 0.3 |
| BSQ | 0.41 | 0.41 | 0.4 | 0.27 | 1 | 0.41 | 0.34 | 0.22 | 0.25 | 0.28 | 0.29 | 0.28 | 0.32 | 0.29 | 0.32 | 0.21 | 0.36 | 0.26 | 0.23 | 0.37 | 0.44 | 0.33 | 0.38 | 0.3 | 0.32 | 0.26 | 0.28 | 0.42 | 0.25 | 0.32 | 0.34 | 0.25 | 0.27 | 0.3 | 0.34 | 0.36 | 0.3 |
| CAK | 0.57 | 0.45 | 0.63 | 0.24 | 0.41 | 1 | 0.55 | 0.21 | 0.39 | 0.4 | 0.4 | 0.38 | 0.42 | 0.3 | 0.36 | 0.34 | 0.44 | 0.31 | 0.32 | 0.37 | 0.56 | 0.35 | 0.61 | 0.31 | 0.43 | 0.21 | 0.37 | 0.48 | 0.27 | 0.26 | 0.41 | 0.29 | 0.38 | 0.39 | 0.29 | 0.41 | 0.3 |
| CHA | 0.6 | 0.41 | 0.61 | 0.27 | 0.34 | 0.55 | 1 | 0.24 | 0.41 | 0.41 | 0.6 | 0.41 | 0.46 | 0.33 | 0.43 | 0.39 | 0.43 | 0.43 | 0.39 | 0.43 | 0.47 | 0.33 | 0.59 | 0.39 | 0.45 | 0.24 | 0.5 | 0.46 | 0.3 | 0.35 | 0.57 | 0.33 | 0.43 | 0.46 | 0.33 | 0.37 | 0.3 |
| CHI | 0.22 | 0.2 | 0.26 | 0.25 | 0.22 | 0.21 | 0.24 | 1 | 0.16 | 0.2 | 0.27 | 0.37 | 0.22 | 0.28 | 0.26 | 0.19 | 0.19 | 0.31 | 0.28 | 0.18 | 0.25 | 0.18 | 0.23 | 0.23 | 0.25 | 0.29 | 0.3 | 0.2 | 0.33 | 0.25 | 0.3 | 0.29 | 0.33 | 0.33 | 0.25 | 0.25 | 0.2 |
| CHO | 0.41 | 0.32 | 0.37 | 0.18 | 0.25 | 0.39 | 0.41 | 0.16 | 1 | 0.29 | 0.28 | 0.24 | 0.31 | 0.21 | 0.26 | 0.22 | 0.27 | 0.24 | 0.22 | 0.27 | 0.33 | 0.27 | 0.45 | 0.23 | 0.27 | 0.17 | 0.26 | 0.32 | 0.21 | 0.26 | 0.31 | 0.21 | 0.24 | 0.27 | 0.2 | 0.29 | 0.3 |
| EFI | 0.51 | 0.34 | 0.4 | 0.22 | 0.28 | 0.4 | 0.41 | 0.2 | 0.29 | 1 | 0.32 | 0.27 | 0.36 | 0.23 | 0.31 | 0.27 | 0.32 | 0.26 | 0.23 | 0.32 | 0.33 | 0.32 | 0.51 | 0.26 | 0.35 | 0.18 | 0.29 | 0.37 | 0.22 | 0.27 | 0.41 | 0.2 | 0.29 | 0.27 | 0.25 | 0.32 | 0.3 |
| ENG | 0.55 | 0.38 | 0.54 | 0.32 | 0.29 | 0.4 | 0.6 | 0.27 | 0.28 | 0.32 | 1 | 0.58 | 0.47 | 0.47 | 0.48 | 0.37 | 0.32 | 0.52 | 0.5 | 0.34 | 0.39 | 0.25 | 0.5 | 0.46 | 0.54 | 0.33 | 0.6 | 0.37 | 0.44 | 0.35 | 0.58 | 0.5 | 0.56 | 0.6 | 0.41 | 0.31 | 0.3 |
| EST | 0.46 | 0.3 | 0.46 | 0.36 | 0.28 | 0.38 | 0.41 | 0.37 | 0.24 | 0.27 | 0.58 | 1 | 0.46 | 0.49 | 0.48 | 0.3 | 0.28 | 0.58 | 0.49 | 0.25 | 0.39 | 0.21 | 0.41 | 0.38 | 0.48 | 0.39 | 0.63 | 0.29 | 0.46 | 0.31 | 0.58 | 0.6 | 0.8 | 0.79 | 0.42 | 0.31 | 0.3 |
| EWE | 0.54 | 0.33 | 0.47 | 0.27 | 0.32 | 0.42 | 0.46 | 0.22 | 0.31 | 0.36 | 0.47 | 0.46 | 1 | 0.29 | 0.42 | 0.29 | 0.35 | 0.4 | 0.33 | 0.3 | 0.38 | 0.31 | 0.52 | 0.31 | 0.52 | 0.25 | 0.46 | 0.32 | 0.32 | 0.31 | 0.49 | 0.35 | 0.43 | 0.49 | 0.32 | 0.33 | 0.3 |
| FIJ | 0.35 | 0.26 | 0.34 | 0.31 | 0.29 | 0.3 | 0.33 | 0.28 | 0.21 | 0.23 | 0.47 | 0.49 | 0.29 | 1 | 0.32 | 0.19 | 0.22 | 0.36 | 0.47 | 0.22 | 0.34 | 0.21 | 0.36 | 0.73 | 0.34 | 0.36 | 0.43 | 0.28 | 0.35 | 0.32 | 0.38 | 0.66 | 0.45 | 0.48 | 0.67 | 0.33 | 0. |
| FIN | 0.46 | 0.31 | 0.42 | 0.33 | 0.32 | 0.36 | 0.43 | 0.26 | 0.26 | 0.31 | 0.48 | 0.48 | 0.42 | 0.32 | 1 | 0.42 | 0.31 | 0.45 | 0.36 | 0.33 | 0.38 | 0.24 | 0.45 | 0.3 | 0.5 | 0.28 | 0.48 | 0.36 | 0.37 | 0.28 | 0.48 | 0.38 | 0.46 | 0.48 | 0.3 | 0.27 | 0.3 |
| FIS | 0.35 | 0.29 | 0.4 | 0.21 | 0.21 | 0.34 | 0.39 | 0.19 | 0.22 | 0.27 | 0.37 | 0.3 | 0.29 | 0.19 | 0.42 | 1 | 0.32 | 0.36 | 0.25 | 0.29 | 0.35 | 0.21 | 0.39 | 0.22 | 0.41 | 0.21 | 0.34 | 0.32 | 0.26 | 0.19 | 0.34 | 0.25 | 0.32 | 0.33 | 0.2 | 0.19 | 0.2 |
| FRE | 0.45 | 0.4 | 0.41 | 0.26 | 0.36 | 0.44 | 0.43 | 0.19 | 0.27 | 0.32 | 0.32 | 0.28 | 0.35 | 0.22 | 0.31 | 0.32 | 1 | 0.3 | 0.22 | 0.43 | 0.43 | 0.52 | 0.39 | 0.24 | 0.36 | 0.21 | 0.35 | 0.51 | 0.25 | 0.26 | 0.34 | 0.22 | 0.26 | 0.28 | 0.23 | 0.35 | 0.3 |
| GEB | 0.44 | 0.33 | 0.44 | 0.32 | 0.26 | 0.31 | 0.43 | 0.31 | 0.24 | 0.26 | 0.52 | 0.58 | 0.4 | 0.36 | 0.45 | 0.36 | 0.3 | 1 | 0.41 | 0.26 | 0.36 | 0.2 | 0.41 | 0.36 | 0.49 | 0.31 | 0.57 | 0.3 | 0.39 | 0.25 | 0.46 | 0.48 | 0.56 | 0.59 | 0.33 | 0.22 | 0. |
| GRC | 0.4 | 0.25 | 0.4 | 0.24 | 0.23 | 0.32 | 0.39 | 0.28 | 0.22 | 0.23 | 0.5 | 0.49 | 0.33 | 0.47 | 0.36 | 0.25 | 0.22 | 0.41 | 1 | 0.21 | 0.29 | 0.2 | 0.38 | 0.45 | 0.37 | 0.34 | 0.44 | 0.25 | 0.37 | 0.25 | 0.42 | 0.71 | 0.52 | 0.51 | 0.4 | 0.22 | 0.2 |
| GRE | 0.39 | 0.41 | 0.39 | 0.25 | 0.37 | 0.37 | 0.43 | 0.18 | 0.27 | 0.32 | 0.34 | 0.25 | 0.3 | 0.22 | 0.33 | 0.29 | 0.43 | 0.26 | 0.21 | 1 | 0.38 | 0.38 | 0.42 | 0.24 | 0.31 | 0.23 | 0.28 | 0.42 | 0.24 | 0.27 | 0.39 | 0.2 | 0.24 | 0.28 | 0.24 | 0.35 | 0.3 |
| GUA | 0.52 | 0.53 | 0.57 | 0.31 | 0.44 | 0.56 | 0.47 | 0.25 | 0.33 | 0.33 | 0.39 | 0.39 | 0.38 | 0.34 | 0.38 | 0.35 | 0.43 | 0.36 | 0.29 | 0.38 | 1 | 0.38 | 0.51 | 0.33 | 0.48 | 0.28 | 0.38 | 0.47 | 0.31 | 0.3 | 0.46 | 0.3 | 0.37 | 0.42 | 0.35 | 0.4 | 0. |
| HAI | 0.41 | 0.35 | 0.33 | 0.23 | 0.33 | 0.35 | 0.33 | 0.18 | 0.27 | 0.32 | 0.25 | 0.21 | 0.31 | 0.21 | 0.24 | 0.21 | 0.52 | 0.2 | 0.2 | 0.38 | 0.38 | 1 | 0.39 | 0.23 | 0.25 | 0.21 | 0.22 | 0.44 | 0.21 | 0.25 | 0.28 | 0.18 | 0.19 | 0.23 | 0.21 | 0.42 | 0.3 |
| HAU | 0.72 | 0.49 | 0.67 | 0.26 | 0.38 | 0.61 | 0.59 | 0.23 | 0.45 | 0.51 | 0.5 | 0.41 | 0.52 | 0.36 | 0.45 | 0.39 | 0.39 | 0.41 | 0.38 | 0.42 | 0.51 | 0.39 | 1 | 0.38 | 0.51 | 0.24 | 0.43 | 0.47 | 0.33 | 0.31 | 0.51 | 0.33 | 0.43 | 0.44 | 0.34 | 0.44 | 0.4 |
| HAW | 0.41 | 0.29 | 0.36 | 0.28 | 0.3 | 0.31 | 0.39 | 0.23 | 0.23 | 0.26 | 0.46 | 0.38 | 0.31 | 0.73 | 0.3 | 0.22 | 0.24 | 0.36 | 0.45 | 0.24 | 0.33 | 0.23 | 0.38 | 1 | 0.37 | 0.35 | 0.37 | 0.31 | 0.35 | 0.34 | 0.38 | 0.59 | 0.38 | 0.41 | 0.67 | 0.32 | 0.3 |
| HIN | 0.55 | 0.34 | 0.53 | 0.32 | 0.32 | 0.43 | 0.45 | 0.25 | 0.27 | 0.35 | 0.54 | 0.48 | 0.52 | 0.34 | 0.52 | 0.5 | 0.41 | 0.36 | 0.49 | 0.37 | 0.31 | 0.48 | 0.25 | 0.51 | 0.37 | 1 | 0.31 | 0.51 | 0.35 | 0.35 | 0.33 | 0.57 | 0.41 | 0.53 | 0.53 | 0.37 | 0.3 |
| HUN | 0.23 | 0.28 | 0.24 | 0.25 | 0.26 | 0.21 | 0.24 | 0.29 | 0.17 | 0.18 | 0.33 | 0.39 | 0.25 | 0.36 | 0.28 | 0.21 | 0.21 | 0.31 | 0.34 | 0.23 | 0.28 | 0.21 | 0.24 | 0.35 | 0.31 | 1 | 0.3 | 0.26 | 0.39 | 0.28 | 0.3 | 0.39 | 0.35 | 0.42 | 0.35 | 0.28 | 0.2 |
| ICE | 0.48 | 0.3 | 0.46 | 0.34 | 0.28 | 0.37 | 0.5 | 0.3 | 0.26 | 0.29 | 0.6 | 0.63 | 0.46 | 0.43 | 0.48 | 0.34 | 0.35 | 0.57 | 0.44 | 0.28 | 0.38 | 0.22 | 0.43 | 0.37 | 0.51 | 0.3 | 1 | 0.34 | 0.39 | 0.28 | 0.55 | 0.48 | 0.58 | 0.61 | 0.35 | 0.26 | 0.3 |
| ITA | 0.47 | 0.5 | 0.49 | 0.26 | 0.42 | 0.48 | 0.46 | 0.2 | 0.32 | 0.37 | 0.37 | 0.29 | 0.32 | 0.28 | 0.36 | 0.32 | 0.51 | 0.3 | 0.25 | 0.42 | 0.47 | 0.44 | 0.47 | 0.31 | 0.35 | 0.26 | 0.34 | 1 | 0.24 | 0.36 | 0.26 | 0.27 | 0.33 | 0.28 | 0.44 | 0.3 | |
| JAM | 0.32 | 0.3 | 0.32 | 0.28 | 0.25 | 0.27 | 0.3 | 0.33 | 0.21 | 0.22 | 0.44 | 0.46 | 0.32 | 0.35 | 0.37 | 0.26 | 0.39 | 0.37 | 0.26 | 0.39 | 0.37 | 0.24 | 0.31 | 0.21 | 0.33 | 0.35 | 0.35 | 0.39 | 0.39 | 0.24 | 1 | 0.32 | 0.41 | 0.42 | 0.4 | 0.38 | 0.34 |
| KAT | 0.36 | 0.28 | 0.31 | 0.25 | 0.32 | 0.26 | 0.35 | 0.25 | 0.26 | 0.27 | 0.35 | 0.31 | 0.31 | 0.32 | 0.28 | 0.19 | 0.26 | 0.25 | 0.25 | 0.27 | 0.3 | 0.25 | 0.31 | 0.34 | 0.33 | 0.28 | 0.28 | 0.27 | 0.32 | 1 | 0.38 | 0.27 | 0.3 | 0.33 | 0.35 | 0.34 | 0.2 |
| KHA | 0.55 | 0.37 | 0.49 | 0.32 | 0.34 | 0.41 | 0.57 | 0.3 | 0.31 | 0.41 | 0.58 | 0.58 | 0.49 | 0.38 | 0.48 | 0.34 | 0.34 | 0.46 | 0.42 | 0.39 | 0.46 | 0.28 | 0.51 | 0.38 | 0.57 | 0.3 | 0.55 | 0.36 | 0.41 | 0.38 | 1 | 0.37 | 0.56 | 0.53 | 0.41 | 0.37 | 0.3 |
| LIT | 0.37 | 0.26 | 0.34 | 0.29 | 0.25 | 0.29 | 0.33 | 0.29 | 0.21 | 0.2 | 0.5 | 0.6 | 0.35 | 0.66 | 0.38 | 0.25 | 0.22 | 0.48 | 0.71 | 0.2 | 0.3 | 0.18 | 0.33 | 0.59 | 0.41 | 0.39 | 0.48 | 0.26 | 0.42 | 0.27 | 0.37 | 1 | 0.56 | 0.68 | 0.55 | 0.25 | 0.2 |
| LIV | 0.47 | 0.28 | 0.46 | 0.33 | 0.27 | 0.38 | 0.43 | 0.33 | 0.24 | 0.29 | 0.56 | 0.8 | 0.43 | 0.45 | 0.46 | 0.32 | 0.26 | 0.56 | 0.52 | 0.24 | 0.37 | 0.19 | 0.43 | 0.38 | 0.53 | 0.35 | 0.58 | 0.27 | 0.4 | 0.3 | 0.56 | 0.56 | 1 | 0.75 | 0.39 | 0.26 | 0.3 |
| LTV | 0.45 | 0.33 | 0.49 | 0.34 | 0.3 | 0.39 | 0.46 | 0.33 | 0.27 | 0.27 | 0.6 | 0.79 | 0.49 | 0.48 | 0.48 | 0.33 | 0.28 | 0.59 | 0.51 | 0.28 | 0.42 | 0.23 | 0.44 | 0.41 | 0.53 | 0.42 | 0.61 | 0.33 | 0.38 | 0.33 | 0.53 | 0.68 | 0.75 | 1 | 0.43 | 0.31 | 0.3 |
| MAO | 0.36 | 0.28 | 0.34 | 0.3 | 0.34 | 0.29 | 0.33 | 0.25 | 0.2 | 0.25 | 0.41 | 0.42 | 0.32 | 0.67 | 0.3 | 0.2 | 0.23 | 0.33 | 0.4 | 0.24 | 0.35 | 0.21 | 0.34 | 0.67 | 0.37 | 0.35 | 0.35 | 0.28 | 0.34 | 0.35 | 0.41 | 0.55 | 0.39 | 0.43 | 1 | 0.3 | 0.3 |
| MAP | 0.43 | 0.4 | 0.37 | 0.27 | 0.36 | 0.41 | 0.37 | 0.25 | 0.29 | 0.32 | 0.31 | 0.31 | 0.33 | 0.33 | 0.27 | 0.19 | 0.35 | 0.22 | 0.22 | 0.35 | 0.49 | 0.42 | 0.44 | 0.32 | 0.29 | 0.28 | 0.26 | 0.44 | 0.32 | 0.34 | 0.37 | 0.25 | 0.26 | 0.31 | 0.3 | 1 | 0.4 |
| MAR | 0.42 | 0.4 | 0.39 | 0.32 | 0.36 | 0.39 | 0.38 | 0.24 | 0.32 | 0.32 | 0.35 | 0.33 | 0.34 | 0.3 | 0.3 | 0.35 | 0.26 | 0.35 | 0.3 | 0.25 | 0.31 | 0.46 | 0.33 | 0.44 | 0.31 | 0.38 | 0.28 | 0.32 | 0.38 | 0.29 | 0.31 | 0.37 | 0.26 | 0.31 | 0.35 | 0.32 | 0.46 |
| MIS | 0.64 | 0.39 | 0.56 | 0.29 | 0.36 | 0.5 | 0.6 | 0.23 | 0.36 | 0.43 | 0.6 | 0.41 | 0.5 | 0.36 | 0.44 | 0.34 | 0.39 | 0.41 | 0.38 | 0.39 | 0.48 | 0.36 | 0.67 | 0.38 | 0.5 | 0.25 | 0.46 | 0.45 | 0.3 | 0.32 | 0.54 | 0.33 | 0.44 | 0.47 | 0.34 | 0.43 | 0.4 |
| MIZ | 0.48 | 0.34 | 0.42 | 0.27 | 0.34 | 0.4 | 0.39 | 0.22 | 0.29 | 0.31 | 0.37 | 0.35 | 0.33 | 0.57 | 0.28 | 0.2 | 0.3 | 0.3 | 0.37 | 0.3 | 0.4 | 0.29 | 0.47 | 0.6 | 0.35 | 0.3 | 0.32 | 0.34 | 0.3 | 0.32 | 0.46 | 0.43 | 0.33 | 0.34 | 0.59 | 0.36 | 0.3 |
| MLT | 0.64 | 0.58 | 0.72 | 0.29 | 0.42 | 0.57 | 0.54 | 0.24 | 0.33 | 0.43 | 0.52 | 0.41 | 0.46 | 0.32 | 0.42 | 0.35 | 0.43 | 0.41 | 0.34 | 0.37 | 0.62 | 0.36 | 0.61 | 0.34 | 0.48 | 0.25 | 0.45 | 0.52 | 0.31 | 0.28 | 0.52 | 0.31 | 0.42 | 0.48 | 0.33 | 0.39 | 0.4 |
| MRD | 0.28 | 0.36 | 0.37 | 0.29 | 0.33 | 0.46 | 0.5 | 0.23 | 0.33 | 0.33 | 0.38 | 0.32 | 0.37 | 0.28 | 0.39 | 0.35 | 0.4 | 0.36 | 0.31 | 0.32 | 0.45 | 0.35 | 0.45 | 0.31 | 0.42 | 0.22 | 0.39 | 0.41 | 0.27 | 0.27 | 0.4 | 0.28 | 0.4 | 0.43 | 0.33 | 0. | |
| MRS | 0.51 | 0.32 | 0.45 | 0.3 | 0.33 | 0.38 | 0.5 | 0.27 | 0.29 | 0.34 | 0.52 | 0.44 | 0.36 | 0.67 | 0.37 | 0.27 | 0.29 | 0.41 | 0.55 | 0.27 | 0.38 | 0.26 | 0.5 | 0.64 | 0.41 | 0.31 | 0.44 | 0.36 | 0.4 | 0.36 | 0.46 | 0.61 | 0.44 | 0.46 | 0.58 | 0.47 | 0.4 |
| NAM | 0.62 | 0.46 | 0.59 | 0.29 | 0.37 | 0.52 | 0.56 | 0.21 | 0.39 | 0.46 | 0.47 | 0.38 | 0.49 | 0.36 | 0.41 | 0.29 | 0.42 | 0.33 | 0.33 | 0.44 | 0.5 | 0.37 | 0.68 | 0.37 | 0.49 | 0.26 | 0.39 | 0.47 | 0.29 | 0.35 | 0.54 | 0.3 | 0.37 | 0.42 | 0.34 | 0.45 | 0.4 |
| NIC | 0.48 | 0.33 | 0.41 | 0.24 | 0.28 | 0.39 | 0.52 | 0.21 | 0.29 | 0.33 | 0.41 | 0.35 | 0.36 | 0.28 | 0.32 | 0.3 | 0.37 | 0.31 | 0.33 | 0.33 | 0.41 | 0.26 | 0.44 | 0.34 | 0.4 | 0.24 | 0.34 | 0.35 | 0.27 | 0.3 | 0.42 | 0.29 | 0.34 | 0.36 | 0.3 | 0.3 | 0.3 |
| OSS | 0.29 | 0.23 | 0.34 | 0.3 | 0.22 | 0.25 | 0.28 | 0.33 | 0.17 | 0.19 | 0.46 | 0.58 | 0.31 | 0.64 | 0.31 | 0.21 | 0.17 | 0.46 | 0.62 | 0.17 | 0.27 | 0.16 | 0.3 | 0.5 | 0.34 | 0.31 | 0.42 | 0.2 | 0.4 | 0.25 | 0.33 | 0.8 | 0.5 | 0.57 | 0.5 | 0.23 | 0.2 |
| PAP | 0.49 | 0.41 | 0.45 | 0.35 | 0.4 | 0.43 | 0.53 | 0.28 | 0.33 | 0.33 | 0.51 | 0.49 | 0.4 | 0.41 | 0.43 | 0.29 | 0.38 | 0.42 | 0.35 | 0.42 | 0.51 | 0.3 | 0.5 | 0.41 | 0.46 | 0.33 | 0.48 | 0.43 | 0.4 | 0.35 | 0.57 | 0.37 | 0.44 | 0.47 | 0.39 | 0.48 | 0.4 |
| PIR | 0.63 | 0.47 | 0.59 | 0.31 | 0.44 | 0.58 | 0.52 | 0.22 | 0.36 | 0.46 | 0.41 | 0.4 | 0.45 | 0.34 | 0.41 | 0.31 | 0.45 | 0.34 | 0.33 | 0.42 | 0.59 | 0.41 | 0.61 | 0.35 | 0.43 | 0.26 | 0.39 | 0.54 | 0.32 | 0.32 | 0.53 | 0.31 | 0.37 | 0.42 | 0.34 | 0.52 | 0.4 |
| PIT | 0.55 | 0.4 | 0.48 | 0.29 | 0.38 | 0.48 | 0.53 | 0.25 | 0.31 | 0.39 | 0.47 | 0.41 | 0.4 | 0.52 | 0.35 | 0.29 | 0.4 | 0.36 | 0.4 | 0.37 | 0.47 | 0.3 | 0.58 | 0.49 | 0.45 | 0.28 | 0.41 | 0.42 | 0.35 | 0.32 | 0.48 | 0.44 | 0.39 | 0.41 | 0.49 | 0.36 | 0.3 |
| RHA | 0.33 | 0.3 | 0.35 | 0.37 | 0.3 | 0.27 | 0.3 | 0.31 | 0.16 | 0.21 | 0.44 | 0.48 | 0.32 | 0.41 | 0.35 | 0.24 | 0.29 | 0.45 | 0.3 | 0.25 | 0.35 | 0.21 | 0.3 | 0.38 | 0.41 | 0.37 | 0.46 | 0.3 | 0.43 | 0.29 | 0.41 | 0.39 | 0.43 | 0.46 | 0.38 | 0.26 | 0.3 |
| RUM | 0.64 | 0.52 | 0.63 | 0.3 | 0.43 | 0.58 | 0.56 | 0.22 | 0.39 | 0.44 | 0.53 | 0.4 | 0.5 | 0.32 | 0.44 | 0.37 | 0.48 | 0.38 | 0.34 | 0.46 | 0.55 | 0.41 | 0.62 | 0.37 | 0.5 | 0.23 | 0.42 | 0.58 | 0.32 | 0.31 | 0.53 | 0.33 | 0.38 | 0.43 | 0.33 | 0.42 | 0.4 |
| RUS | 0.29 | 0.23 | 0.33 | 0.25 | 0.21 | 0.27 | 0.3 | 0.3 | 0.18 | 0.18 | 0.46 | 0.51 | 0.3 | 0.54 | 0.31 | 0.2 | 0.18 | 0.4 | 0.65 | 0.18 | 0.26 | 0.18 | 0.31 | 0.51 | 0.34 | 0.36 | 0.39 | 0.22 | 0.39 | 0.25 | 0.35 | 0.77 | 0.43 | 0.54 | 0.45 | 0.25 | 0.2 |
| SAA | 0.52 | 0.32 | 0.51 | 0.31 | 0.28 | 0.42 | 0.54 | 0.31 | 0.26 | 0.32 | 0.59 | 0.65 | 0.47 | 0.4 | 0.55 | 0.4 | 0.33 | 0.55 | 0.48 | 0.27 | 0.4 | 0.23 | 0.48 | 0.37 | 0.55 | 0.32 | 0.65 | 0.36 | 0.41 | 0.29 | 0.56 | 0.5 | 0.62 | 0.63 | 0.34 | 0.28 | 0.3 |
| SAM | 0.26 | 0.22 | 0.24 | 0.22 | 0.2 | 0.22 | 0.29 | 0.24 | 0.19 | 0.23 | 0.28 | 0.28 | 0.23 | 0.31 | 0.24 | 0.18 | 0.23 | 0.21 | 0.24 | 0.2 | 0.24 | 0.2 | 0.27 | 0.28 | 0.25 | 0.25 | 0.27 | 0.24 | 0.24 | 0.24 | 0.28 | 0.24 | 0.24 | 0.25 | 0.26 | 0.27 | 0.2 |
| SGO | 0.42 | 0.37 | 0.4 | 0.31 | 0.31 | 0.3 | 0.33 | 0.27 | 0.29 | 0.28 | 0.45 | 0.45 | 0.39 | 0.31 | 0.34 | 0.22 | 0.26 | 0.38 | 0.36 | 0.27 | 0.39 | 0.26 | 0.4 | 0.3 | 0.43 | 0.29 | 0.39 | 0.28 | 0.41 | 0.25 | 0.47 | 0.4 | 0.41 | 0.44 | 0.32 | 0.3 | 0.3 |
| SNG | 0.66 | 0.48 | 0.62 | 0.27 | 0.36 | 0.55 | 0.61 | 0.22 | 0.44 | 0.45 | 0.51 | 0.41 | 0.49 | 0.32 | 0.39 | 0.32 | 0.4 | 0.36 | 0.36 | 0.45 | 0.5 | 0.36 | 0.64 | 0.36 | 0.48 | 0.25 | 0.46 | 0.46 | 0.33 | 0.33 | 0.53 | 0.33 | 0.4 | 0.43 | 0.32 | 0.43 | 0.4 |
| SPA | 0.53 | 0.52 | 0.61 | 0.24 | 0.45 | 0.62 | 0.48 | 0.21 | 0.38 | 0.37 | 0.38 | 0.34 | 0.37 | 0.29 | 0.35 | 0.28 | 0.47 | 0.27 | 0.28 | 0.43 | 0.56 | 0.42 | 0.62 | 0.32 | 0.37 | 0.27 | 0.31 | 0.6 | 0.28 | 0.26 | 0.4 | 0.28 | 0.32 | 0.35 | 0.3 | 0.51 | 0.3 |
| SRA | 0.42 | 0.31 | 0.41 | 0.29 | 0.3 | 0.3 | 0.39 | 0.34 | 0.25 | 0.28 | 0.44 | 0.49 | 0.38 | 0.3 | 0.39 | 0.28 | 0.29 | 0.52 | 0.42 | 0.3 | 0.34 | 0.23 | 0.37 | 0.32 | 0.42 | 0.29 | 0.48 | 0.29 | 0.38 | 0.28 | 0.46 | 0.4 | 0.46 | 0.43 | 0.32 | 0.25 | 0.2 |
| SWA | 0.57 | 0.44 | 0.49 | 0.22 | 0.36 | 0.54 | 0.43 | 0.17 | 0.37 | 0.44 | 0.31 | 0.28 | 0.36 | 0.27 | 0.34 | 0.28 | 0.39 | 0.24 | 0.26 | 0.38 | 0.5 | 0.47 | 0.59 | 0.27 | 0.33 | 0.22 | 0.29 | 0.48 | 0.24 | 0.25 | 0.38 | 0.24 | 0.25 | 0.38 | 0.3 | | |
| SWE | 0.45 | 0.25 | 0.4 | 0.32 | 0.27 | 0.31 | 0.43 | 0.28 | 0.26 | 0.25 | 0.52 | 0.51 | 0.41 | 0.41 | 0.5 | 0.3 | 0.29 | 0.6 | 0.43 | 0.24 | 0.33 | 0.19 | 0.39 | 0.39 | 0.52 | 0.33 | 0.62 | 0.29 | 0.4 | 0.3 | 0.47 | 0.53 | 0.49 | 0.54 | 0.36 | 0.27 | 0.2 |
| TAG | 0.27 | 0.23 | 0.22 | 0.15 | 0.18 | 0.22 | 0.29 | 0.14 | 0.22 | 0.3 | 0.19 | 0.14 | 0.21 | 0.12 | 0.2 | 0.19 | 0.29 | 0.17 | 0.13 | 0.28 | 0.21 | 0.3 | 0.26 | 0.15 | 0.18 | 0.13 | 0.18 | 0.28 | 0.12 | 0.21 | 0.22 | 0.11 | 0.16 | 0.17 | 0.13 | 0.26 | 0.2 |
| TAJ | 0.39 | 0.29 | 0.44 | 0.29 | 0.25 | 0.3 | 0.41 | 0.3 | 0.23 | 0.24 | 0.58 | 0.52 | 0.34 | 0.4 | 0.4 | 0.28 | 0.24 | 0.47 | 0.51 | 0.23 | 0.38 | 0.19 | 0.38 | 0.38 | 0.44 | 0.3 | 0.49 | 0.28 | 0.39 | 0.28 | 0.44 | 0.53 | 0.49 | 0.49 | 0.37 | 0.24 | 0.2 |
| TKP | 0.36 | 0.32 | 0.34 | 0.34 | 0.35 | 0.3 | 0.33 | 0.33 | 0.23 | 0.25 | 0.43 | 0.53 | 0.4 | 0.37 | 0.41 | 0.22 | 0.28 | 0.41 | 0.4 | 0.29 | 0.39 | 0.26 | 0.33 | 0.35 | 0.39 | 0.38 | 0.41 | 0.3 | 0.5 | 0.31 | 0.46 | 0.44 | 0.47 | 0.49 | 0.35 | 0.32 | 0.3 |
| TOA | 0.45 | 0.37 | 0.37 | 0.34 | 0.38 | 0.36 | 0.42 | 0.26 | 0.26 | 0.31 | 0.45 | 0.41 | 0.35 | 0.37 | 0.37 | 0.32 | 0.24 | 0.35 | 0.35 | 0.28 | 0.34 | 0.42 | 0.35 | 0.4 | 0.39 | 0.36 | 0.33 | 0.38 | 0.39 | 0.37 | 0.38 | 0.41 | 0.36 | 0.44 | 0.4 | 0.4 | 0.3 |
| TOB | 0.41 | 0.36 | 0.37 | 0.26 | 0.29 | 0.36 | 0.39 | 0.16 | 0.27 | 0.32 | 0.32 | 0.28 | 0.34 | 0.24 | 0.34 | 0.24 | 0.36 | 0.29 | 0.41 | 0.28 | 0.23 | 0.35 | 0.39 | 0.27 | 0.38 | 0.24 | 0.36 | 0.22 | 0.34 | 0.39 | 0.26 | 0.23 | 0.38 | 0.23 | 0.24 | 0.28 | 0.2 |
| TON | 0.48 | 0.29 | 0.4 | 0.3 | 0.27 | 0.35 | 0.47 | 0.24 | 0.29 | 0.3 | 0.6 | 0.48 | 0.38 | 0.46 | 0.38 | 0.32 | 0.31 | 0.36 | 0.37 | 0.31 | 0.37 | 0.23 | 0.48 | 0.4 | 0.43 | 0.28 | 0.44 | 0.33 | 0.35 | 0.33 | 0.5 | 0.37 | 0.42 | 0.45 | 0.35 | 0.34 | 0.3 |
| TRK | 0.48 | 0.5 | 0.46 | 0.34 | 0.38 | 0.42 | 0.43 | 0.24 | 0.28 | 0.34 | 0.42 | 0.37 | 0.41 | 0.29 | 0.4 | 0.32 | 0.41 | 0.4 | 0.28 | 0.37 | 0.5 | 0.35 | 0.48 | 0.31 | 0.42 | 0.29 | 0.38 | 0.43 | 0.31 | 0.28 | 0.39 | 0.29 | 0.35 | 0.4 | 0.31 | 0.37 | 0.4 |
| UDM | 0.57 | 0.46 | 0.54 | 0.3 | 0.34 | 0.49 | 0.5 | 0.24 | 0.35 | 0.39 | 0.42 | 0.35 | 0.41 | 0.31 | 0.37 | 0.33 | 0.4 | 0.33 | 0.32 | 0.38 | 0.5 | 0.37 | 0.57 | 0.34 | 0.46 | 0.23 | 0.39 | 0.47 | 0.27 | 0.28 | 0.43 | 0.31 | 0.35 | 0.38 | 0.3 | 0.38 | 0.3 |
| ULA | 0.46 | 0.33 | 0.4 | 0.27 | 0.32 | 0.37 | 0.37 | 0.24 | 0.28 | 0.34 | 0.36 | 0.38 | 0.34 | 0.49 | 0.3 | 0.22 | 0.27 | 0.31 | 0.42 | 0.29 | 0.41 | 0.3 | 0.46 | 0.5 | 0.35 | 0.3 | 0.34 | 0.37 | 0.29 | 0.33 | 0.4 | 0.47 | 0.39 | 0.44 | 0.5 | 0.38 | 0.3 |
| VEP | 0.44 | 0.32 | 0.45 | 0.35 | 0.32 | 0.38 | 0.47 | 0.29 | 0.27 | 0.3 | 0.49 | 0.49 | 0.42 | 0.34 | 0.76 | 0.42 | 0.31 | 0.47 | 0.39 | 0.32 | 0.38 | 0.25 | 0.43 | 0.34 | 0.48 | 0.27 | 0.46 | 0.32 | 0.35 | 0.26 | 0.48 | 0.41 | 0.48 | 0.51 | 0.3 | 0.24 | 0. |
| VIE | 0.46 | 0.36 | 0.42 | 0.26 | 0.34 | 0.37 | 0.45 | 0.22 | 0.32 | 0.29 | 0.41 | 0.34 | 0.31 | 0.36 | 0.33 | 0.23 | 0.32 | 0.32 | 0.28 | 0.33 | 0.41 | 0.3 | 0.42 | 0.4 | 0.38 | 0.28 | 0.36 | 0.41 | 0.28 | 0.36 | 0.4 | 0.31 | 0.31 | 0.37 | 0.36 | 0.36 | 0.3 |
| YAB | 0.33 | 0.32 | 0.29 | 0.19 | 0.27 | 0.31 | 0.29 | 0.17 | 0.27 | 0.32 | 0.25 | 0.25 | 0.25 | 0.28 | 0.26 | 0.19 | 0.33 | 0.2 | 0.22 | 0.29 | 0.32 | 0.33 | 0.39 | 0.25 | 0.26 | 0.24 | 0.25 | 0.33 | 0.21 | 0.24 | 0.3 | 0.21 | 0.24 | 0.26 | 0.22 | 0.32 | 0.2 |
| ZAP | 0.5 | 0.47 | 0.52 | 0.26 | 0.36 | 0.51 | 0.48 | 0.22 | 0.36 | 0.35 | 0.38 | 0.37 | 0.37 | 0.34 | 0.39 | 0.28 | 0.36 | 0.31 | 0.3 | 0.41 | 0.51 | 0.37 | 0.58 | 0.33 | 0.36 | 0.28 | 0.35 | 0.45 | 0.34 | 0.3 | 0.44 | 0.29 | 0.35 | 0.38 | 0.33 | 0.47 | 0.2 |

Height

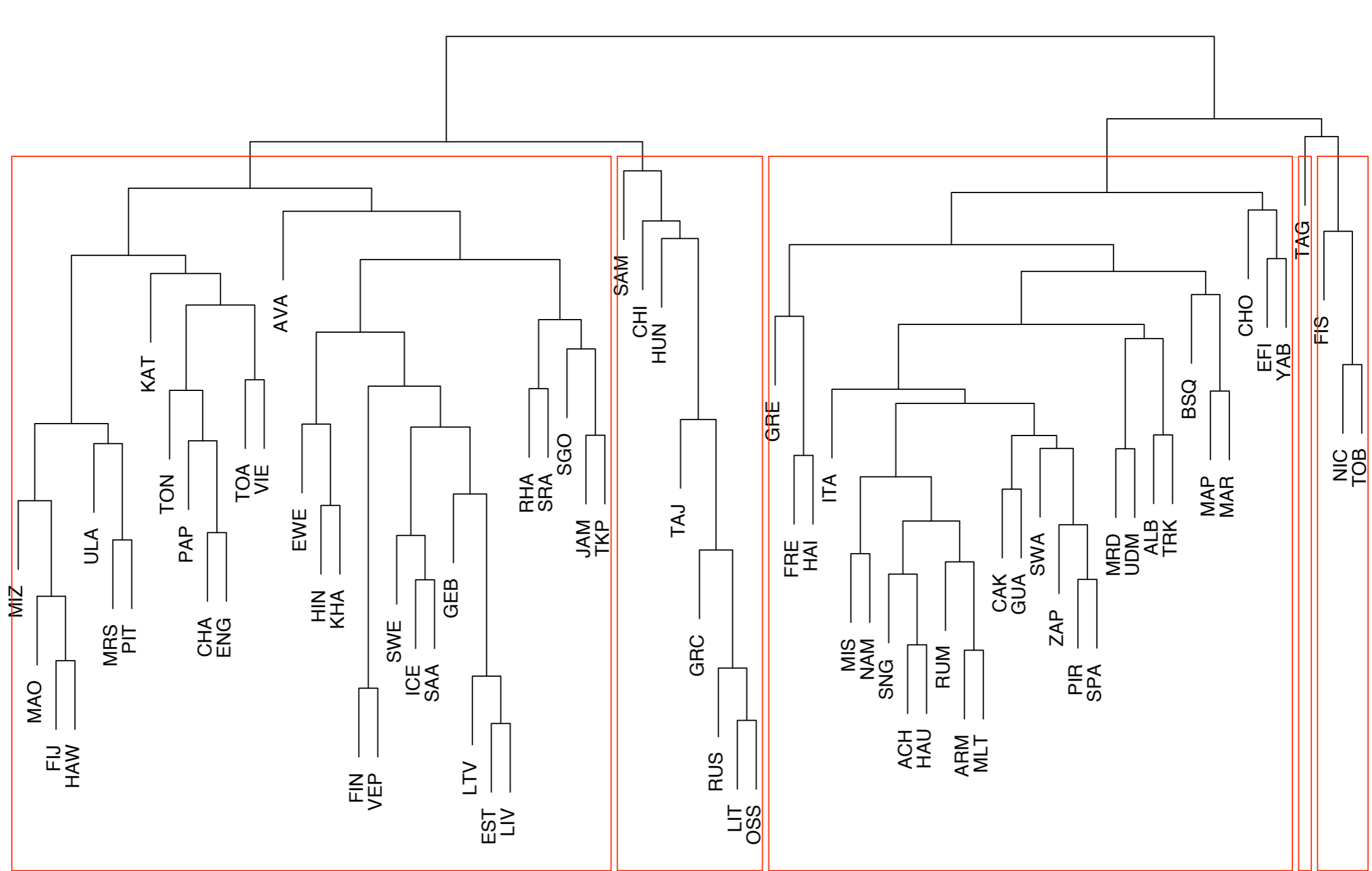hclust (*, "complete")

hclust (*, "complete")

Height

hclust (*, "complete")

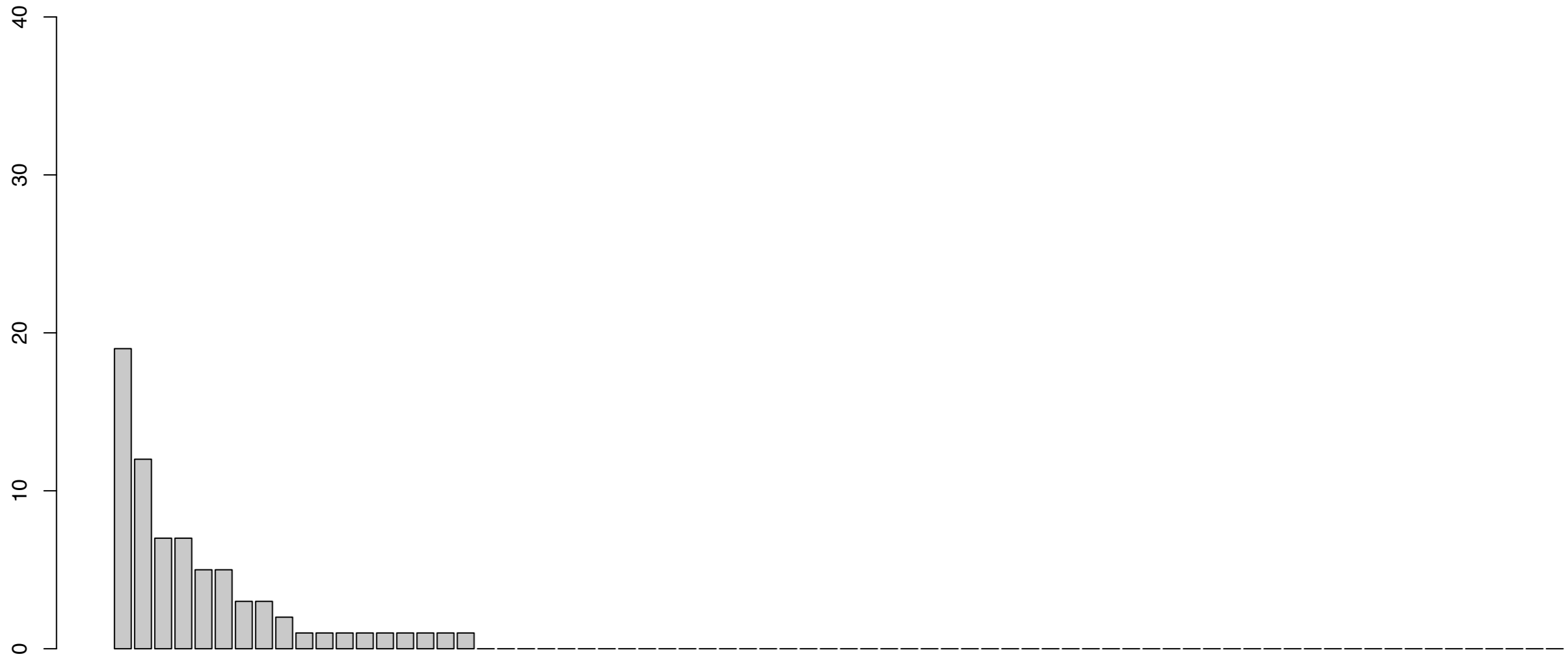Height

hclust (*, "complete")

hclust (*, "complete")

# Where is this going?

- Typological distributions are shaped by different forces:

  ▸ language change

  ▸ choices made by linguists

- Negative exponential distributions for unordered typologies are a promising idea

- The distributions are clearly not even !