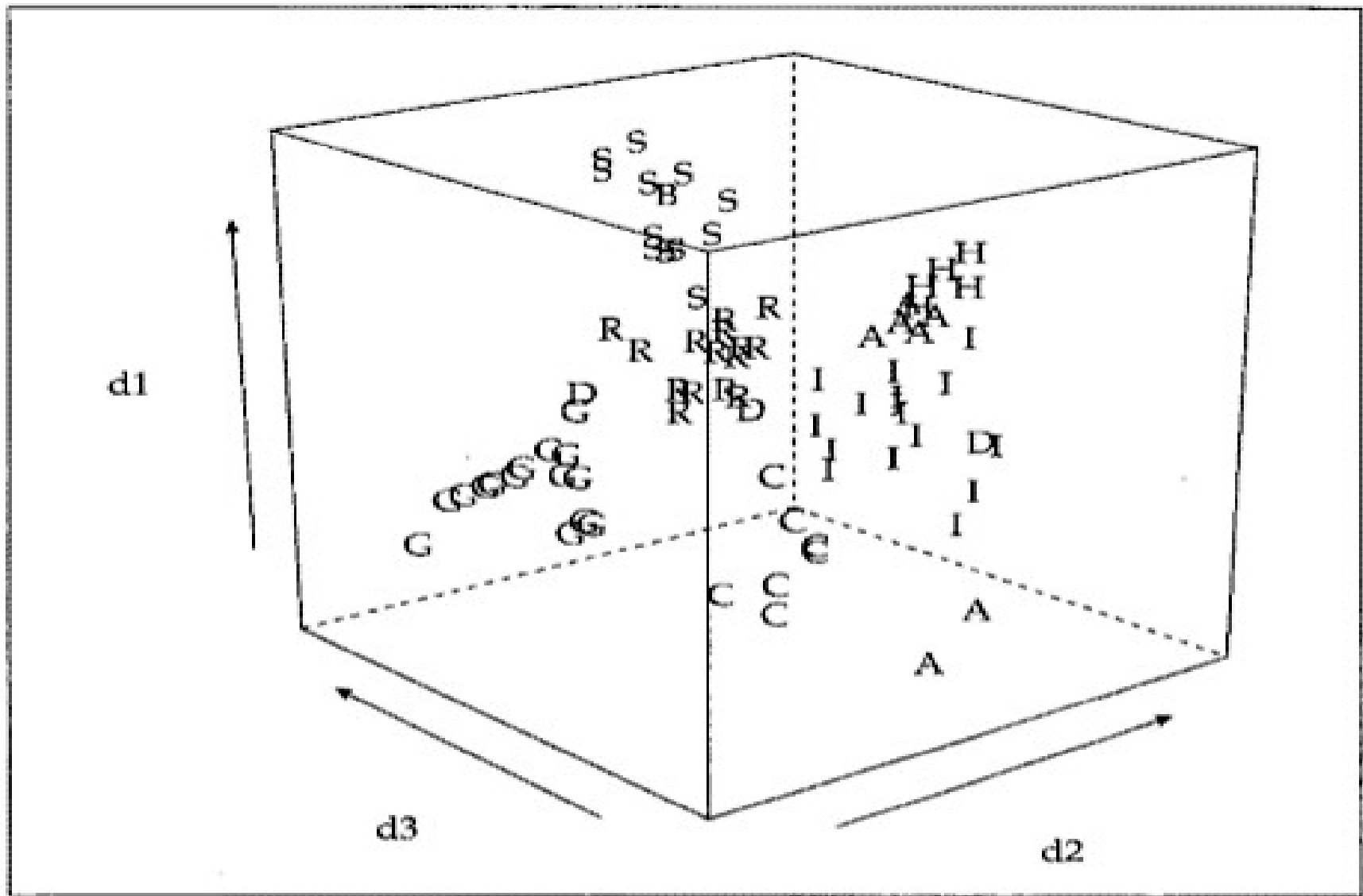


Adatbányászat az R nyelv alkalmazásával

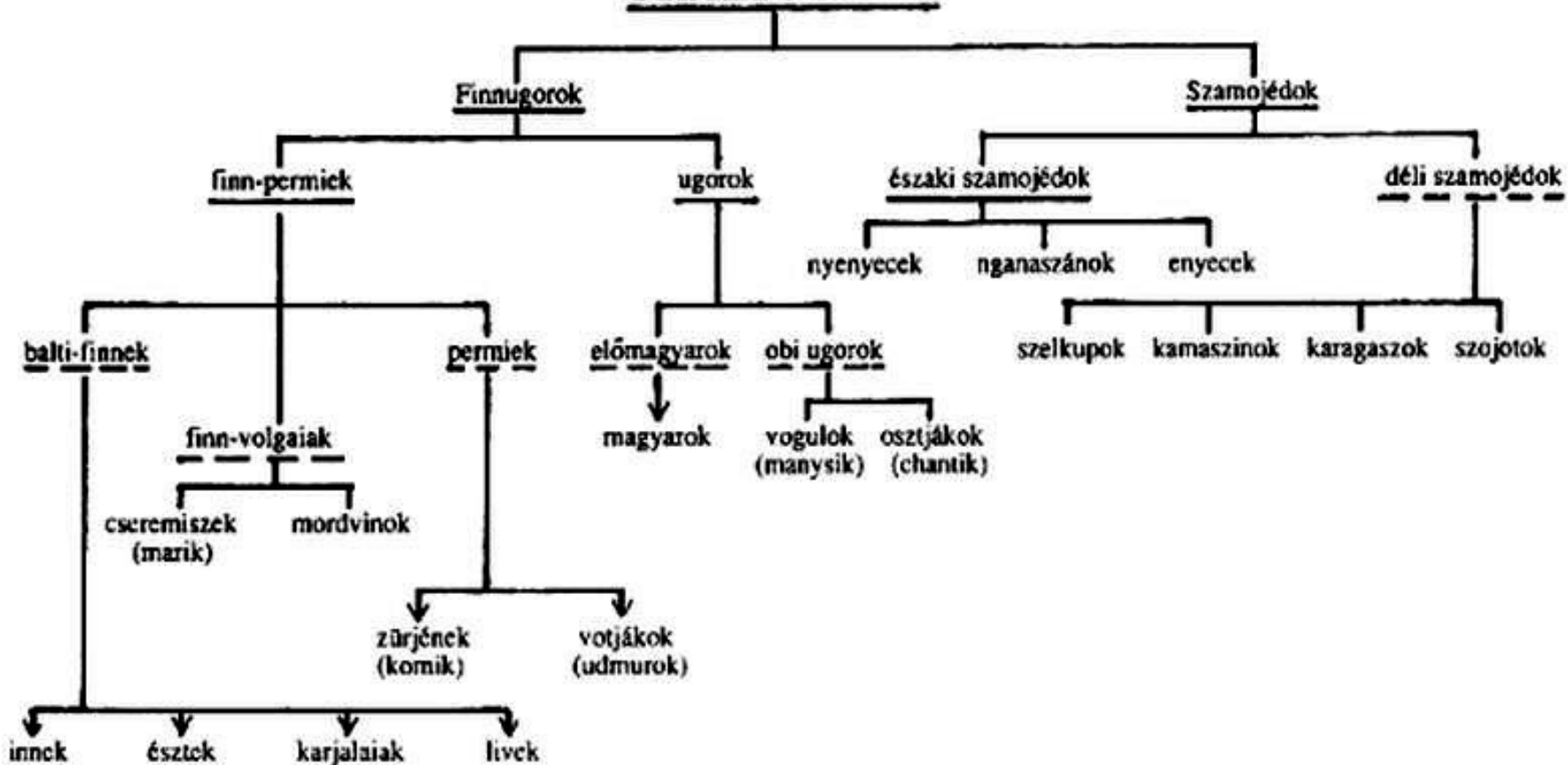
Varjú Zoltán

- Mennyire hasonlítanak egymásra a finnugor nyelvek?
- Milyen más nyelvekre hasonlítanak?
- Hogyan mérhetjük a hasonlóságot?
- Hogyan mutassuk meg a hasonlósági metrikát olyan közönségnek mely tagjai alapvetően nem a számok nyelvén beszélnek?
- Saját gyűjtésekre és online elérhető szövegekre alapozva a karaktorsorok közötti hasonlóságot mértük



A FINNUGOR NÉPEK CSALÁDFÁJA

URALI NÉPEK



- Az R-ben mindenre van megoldás
- Ha hozzá szokik az ember a dokumentáció stílusához, akkor képes igényeire szabni a példa szkripteket
- Learning by doing
- A nyelvtudásnak több szintje van a programozási nyelvek esetében is

- Rövid áttekintés
- R a gyakorlatban
- R és Big Data
- Jövő

- Az R nyelv az S statisztikai programozási nyelv “dialektusa”
- 1976 Bell Labs, John Chambers és tsai Fortranban implementálták az S nyelvet
- 1988 C-ben újraírták az S-t
- 1993-ban Insightful Corp. megkapja az S licencét, ma TIBCO birtokolja S-PLUS néven

- 1991 Ross Ihaka és Robert Gentleman elindítja az R fejlesztését
- 1995 GPL licenc alatt jelenik meg a nyelv
- 2000 R 1.0.0
- 2012.06.22 R 2.15.1 “Roasted Marshmallows”

Kis R történelem és háttér

- CRAN

<http://cran.r-project.org/>

- 5300 csomag

- CRAN Task Views

- ggplot2

<http://ggplot2.org/>

- Wilkson et al. *The Grammar of Graphics*



Miért R?

- CRAN 5300 csomag
- A statisztikai lingua franca
 - A legjobban dokumentált nyelv a területen
- Aktív fejlesztői közösség
- Aktív felhasználói közösség
- Trendi

“[W]e wanted users to be able to begin in an interactive environment, where they did not consciously think of themselves as programming. Then as their needs became clearer and their sophistication increased, they should be able to slide gradually into programming, when the language and system aspects would become more important.” Chambers

(<http://www.stat.bell-labs.com/S/history.html>)

- R, mint adatfelfedezési eszköz (EDA)



- R, mint “number cruncher”



- Elméletileg az R “képeségei” lefedik az egész elemzési folyamatot
- Gyakorlatban az előfeldolgozás általában más nyelven történik (pl. Python, Perl)
- “interaktív” szkriptek a repl segítségével
- A szkriptek “refaktorálása” riportokhoz
- prototyping

- Szoftverfejlesztés vs. algoritmikus elemzés
- Kelly – Mill: *A Software Chasm*
- Az R a legelterjedtebb nyelv kutatói és ipari körökben is
- Prabhu – Jablin: *A survey of practice of computational science*

- In-memory eszköz
- a rendelkezésre álló RAM 10-20 százalékát meghaladó adatmennyiség már lassuláshoz vezet, 50% felett pedig tkp. használhatatlan
- Bigmemory - köztes megoldást kínál desktop környezetben, az ún. memory-mapping technikával egyfajta virtuális memóriát használ
- Snow, Multicore, Parallel, R+Hadoop, RHIFE, Segue – open source megoldások, amik jók, de nem triviális R elősimereteket kíván meg alkalmazásuk

scale

Big Data Statistics for R

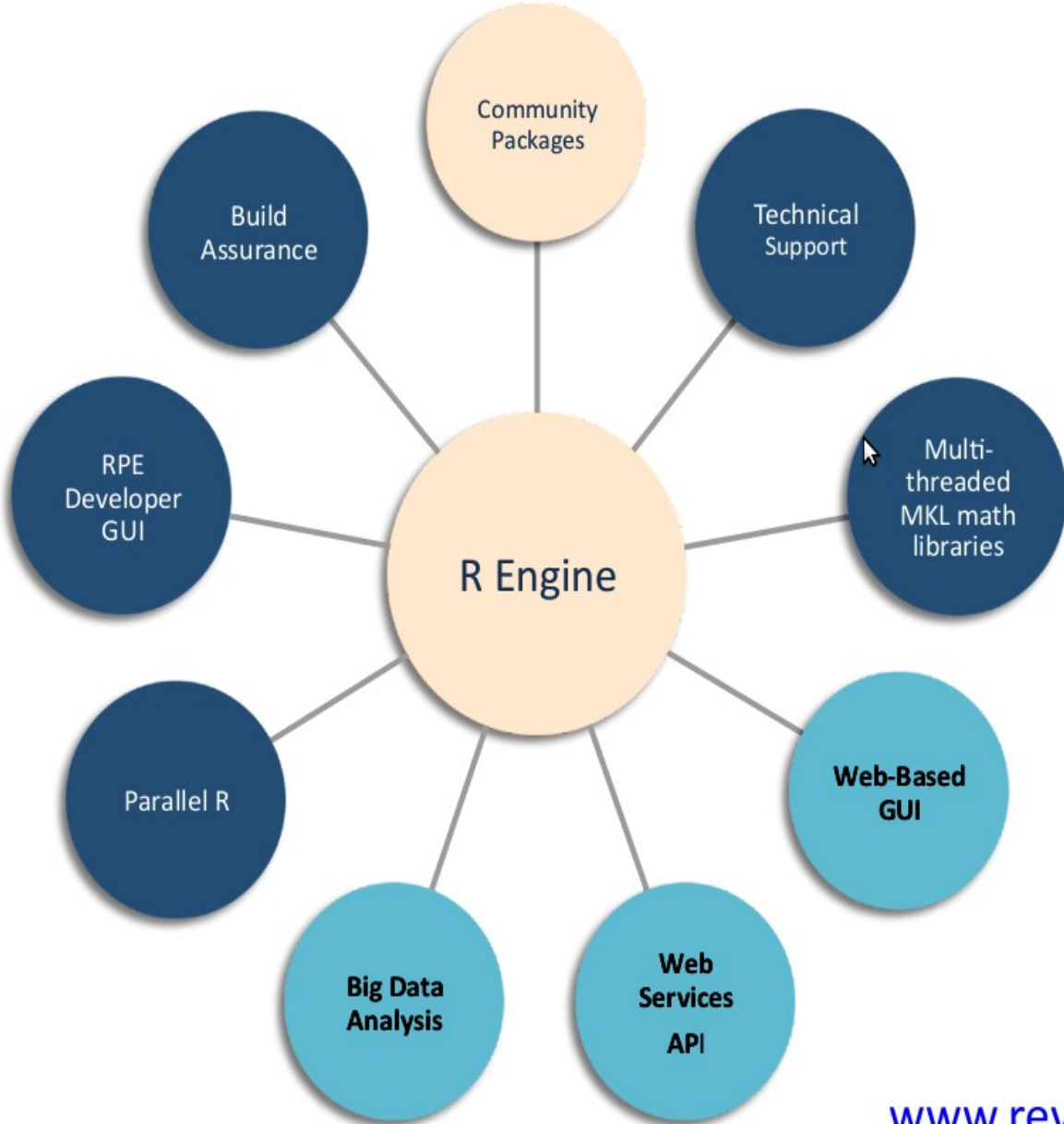
deploy




Application Integration for R



With **Revolution R**
Enterprise

R és Big Data



-  Community - Open Source
-  Revolution – Proprietary additions
-  Revolution – Forthcoming proprietary additions

- Az átlag R felhasználó nem programozó
- Humán faktor: miért tanuljon meg valaki még egy új nyelvet, technológiát stb. ?
- \$1000 / workstation Revolution Analytics licenc díj

- R fejlesztések orvosolják a problémákat
- Egyre többen ismerik a nyelvet
“... *the sexy job in the next 10 years will be statisticians*” (Hal Varian, Google)



cloudnumbers.com

- Funkcionális JVM nyelvek
- Clojure Incanter
- <http://incanter.org/>



Clojure



- Scala breeze
- <https://github.com/scalanlp/breeze>

 **Scala**

Köszönöm a figyelmet!

- Kereső Világ blog

<http://kereses.blog.hu/>

Számítógépes nyelvészet szakmai blog

<http://szamitogepesnyelveszet.blogspot.com/>

- varju.zoltan@weblib.hu

- @zoltanvarju