

# **Social media jako nové pole pro data mining**

30. 5. 2013 Josef Šlerka

# Social data profiling

řekni mi co lajkuješ a já to povím, kdo jsi



ČSSD Hostěradice  
Political Party



OVV ČSSD Pardubice  
Organization



OVV ČSSD Strakonice  
Political Party



Podporuji, aby prezidentem ...  
Community



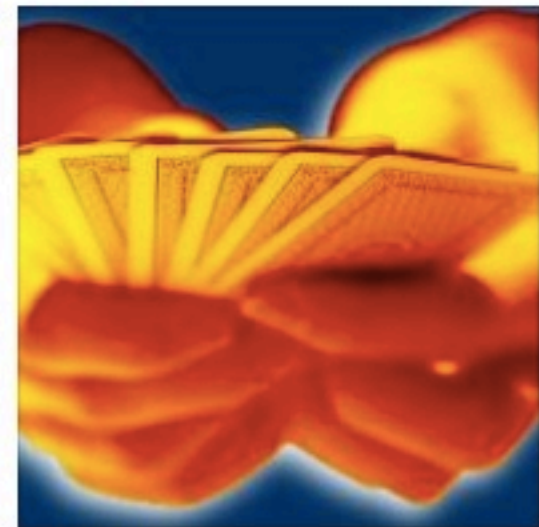
Vladimír Plaček  
Politician



MO ČSSD Krnov  
Political Party



ČSSD Moravská Ostrava a Pří...  
Political Party



Výklad z karet  
Community



Okresní organizace ČSSD Jičín  
Political Party



Vinný sklípek u Kratochvílů  
Restaurant



Deník Legie  
Magazine



Adam Rykala  
Politician





Centre Tchèque České Centrum  
Non-Profit Organization



Kalhoty pro Václava Havla  
Community



Bel Mondo  
Media/News/Publishing



Izrael v České republice  
Consulate & Embassy



Kamil Fila rulez  
Community



Evropská observatoř žurnalis...  
Media/News/Publishing



Woody Allen  
Public Figure



TIME  
Media/News/Publishing



100 let skautingu v českých ...  
Community



causes.com  
App



HNZprávy  
Newspaper



David Lynch  
Artist





Knihobežník  
Games/Toys

UKÁČKO<sup>CZ</sup>  
PORTÁL STUDENTŮ UK

UKáčko.cz – Portál studentů ...  
News/Media Website

FFAKT

FFakt, časopis studentů Filoz...  
Magazine



#procjsemtady  
Community

AKČNÍ  
LETENKY

Akční letenky  
Travel/Leisure



Bloxter  
Local Business



brmlab  
Education



Mládež pražského seniorátu ...  
Community



ČCE Praha 10 – Vršovice  
Religious Organization



qwerly  
Product/Service

CZECH  
O<sup>®</sup>IGINAL  
FASHION

www.czechoriginalfashion.cz  
Website



SNM.gug.cz  
Community

# Prezidentská volba

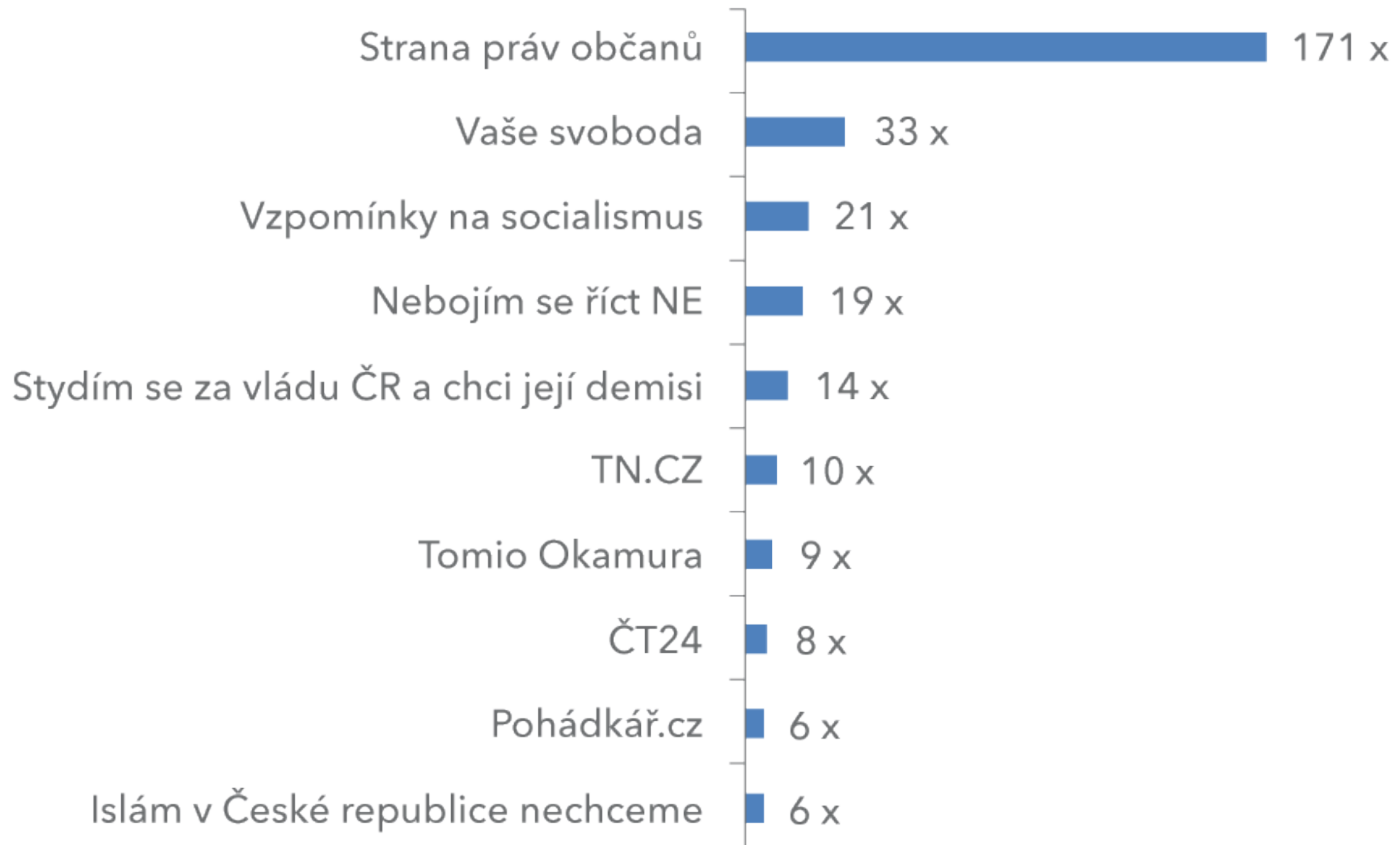
analýza fanoušků prezidentských kandidátů na Facebooku

analyzováno přes 1.600 fanoušků

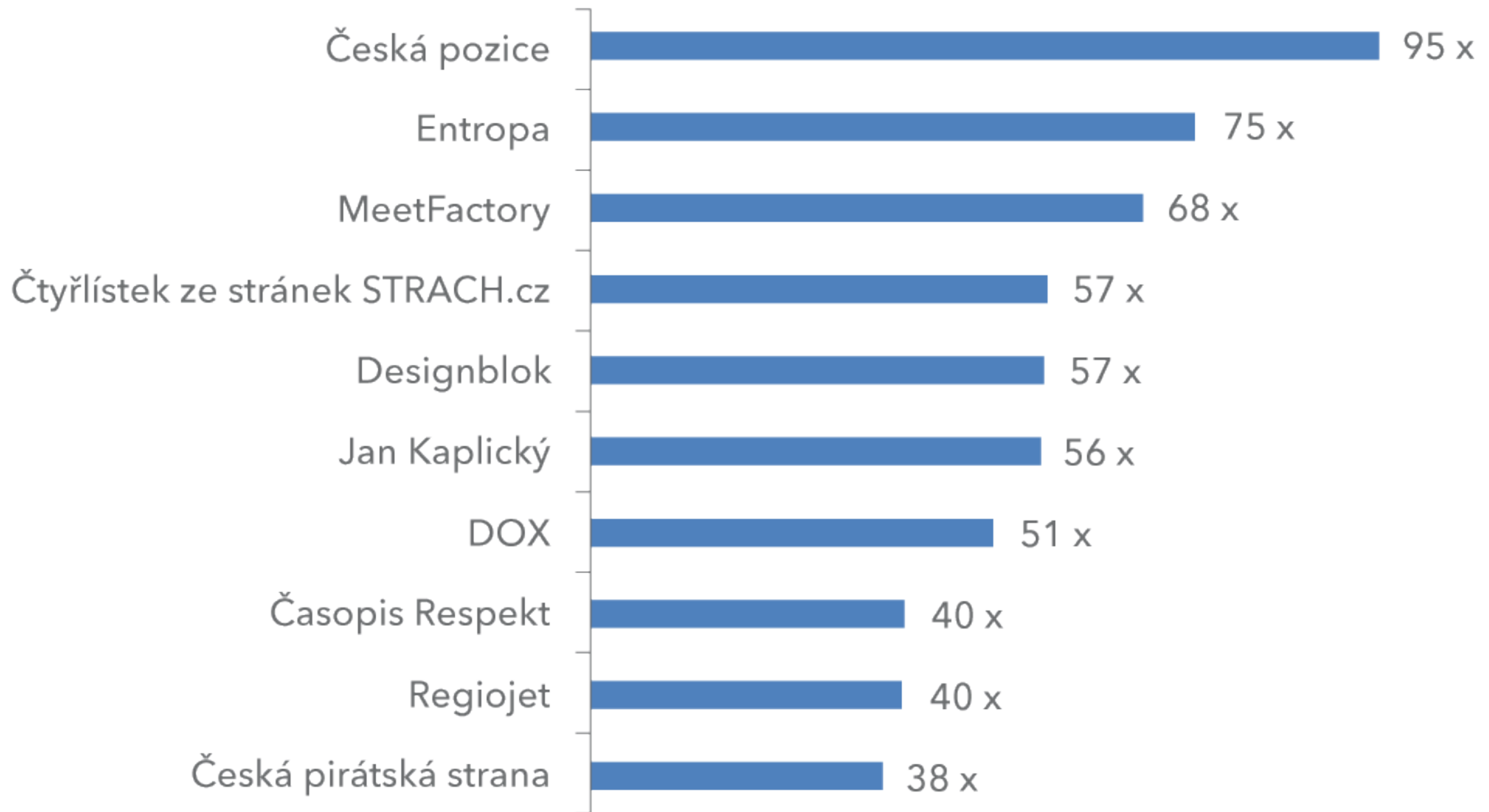
mapa charakteristických odchylek fanoušků od průměrného uživatele

společný projekt České televize a Studia nových médií a blízkých lidí (Josef Šlerka, Jan Schmid)

# Miloš Zeman

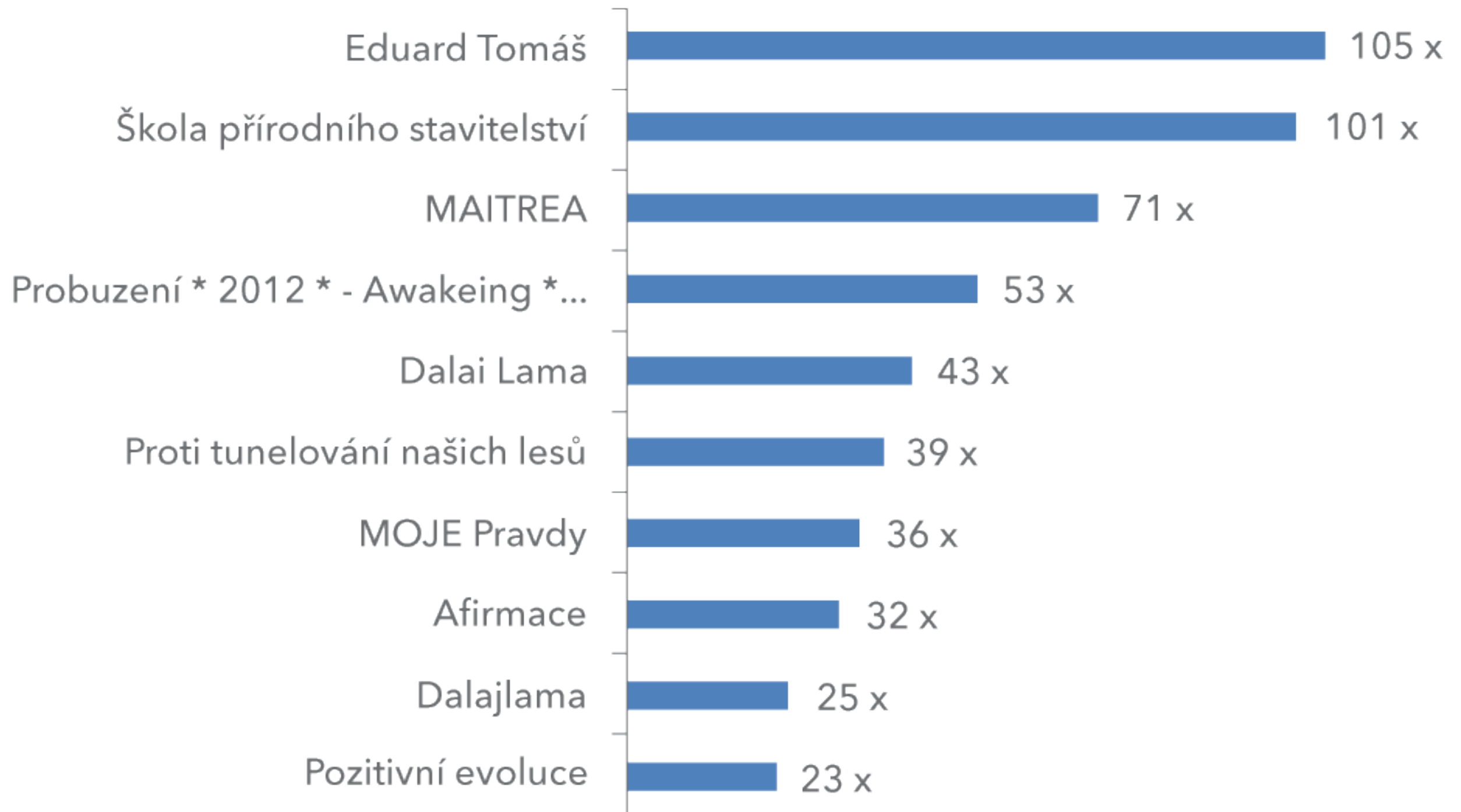


# Vladimír Franz





# Táňa Fišerová



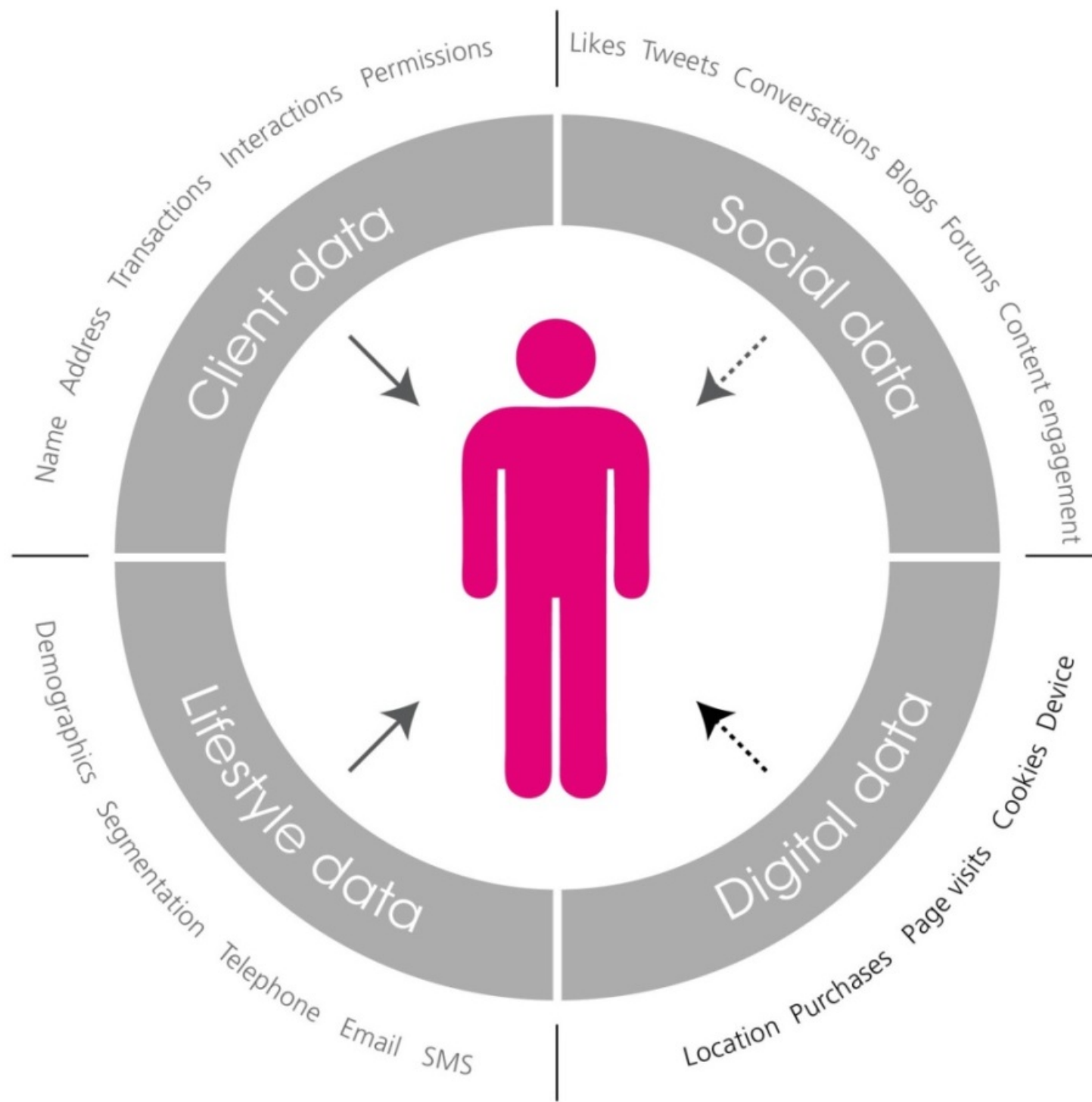
# Pěkný, ale...

... příliš “big picture” pro každodenní business!

# Social Insight Finder







# Základní principy

obchodní segmentace na tvrdých datech (transakční data, webová analytika aj.)

zájmové preference a psycho-demografie na datech ze sociálních sítí (Facebook, Twitter, blogy aj.)

# Případová studie

velký klient z oblasti e-commerce

požadavek na segmentaci klientů a jejich charakteristiku

vhled do segmentu, který utrácí i do segmentu, který neutráčí



# Postup

klasická segmentace z transakčních a dalších dat

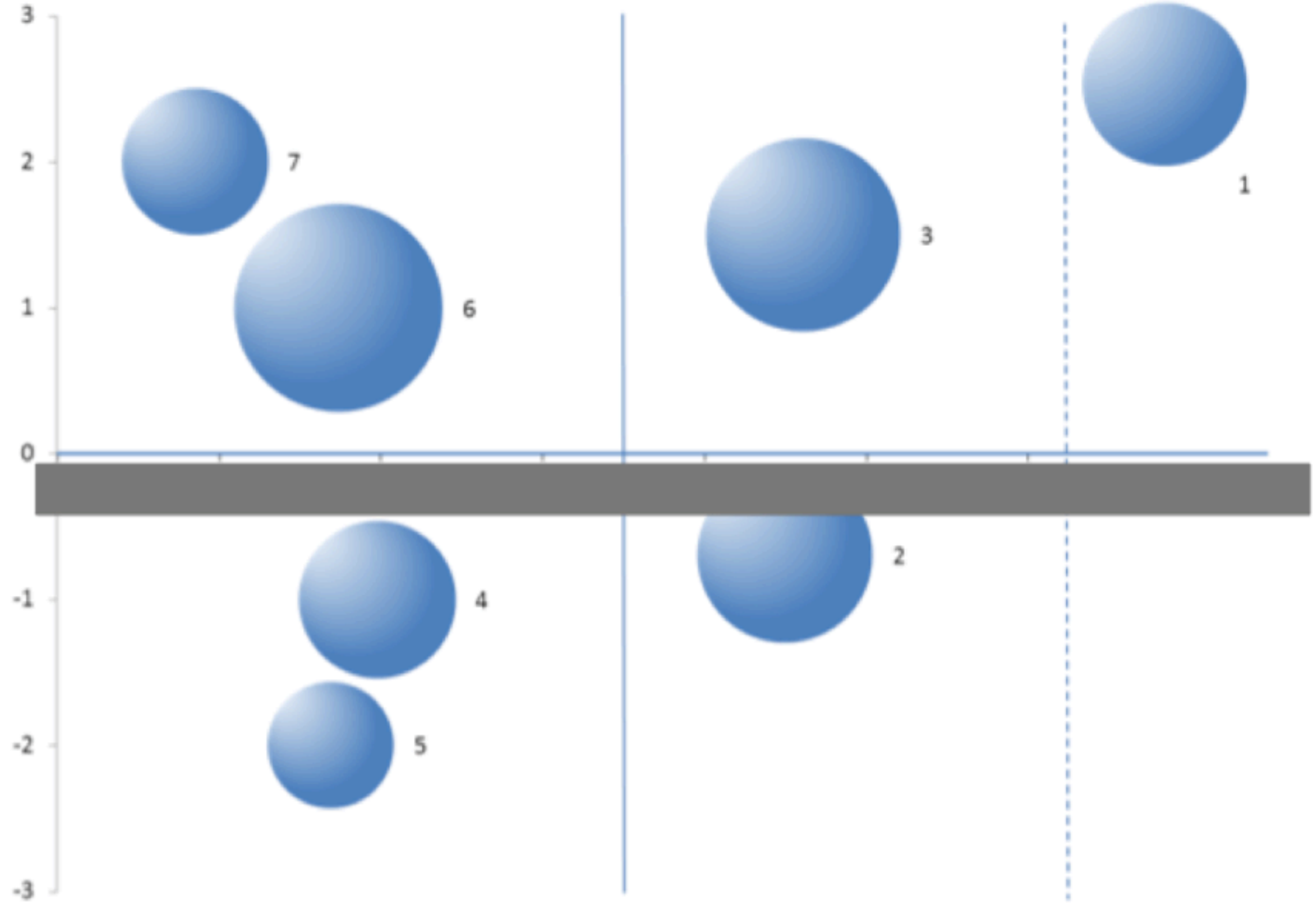
obohacení o typické odchylky v preferencích na Facebooku od průměrné populace v jednotlivých segmentech

propojení obou analýz

**Ukázky z výsledné  
analýzy (trochu)**

# 7 obchodně aktivních segmentů

napětí zábava vs informace



odchylka od mediánu útraty



# Č. 2 URBAN YUPIIES

průměrné revenue

objednávají: goods, food, fun, méně travel

browsing: fun, food

na webu nadprůměrně aktivní ve všech kategoriích

pohlaví: nejvíce mužů z nakupujících segmentů

výrazně nadprůměrně platba kartou

věk: medián XX let



muži: městský starší teenager, po škole, single nebo nemá žádný vážný vztah, má peníze, chce se ukázat

price sensitive: Časopis dTest

technika: MobilMania.cz, Zive.cz

zábava: Fotbálek, Zelená, Show Jana Krause, Maxim

požitkáři: České pivo, JenProMuze.cz, Češky jsou nejkrásnější holky na světě etc.



ženy: studentky po nebo na VŠ (Fakulta roku), single nebo nemá žádný vážný vztah, má peníze pro sebe

chtějí se ukázat, řeší módu, značky – Starbucks, Botyk.cz, Fashion Days, CCC Boty, Zoot, Módní peklo

pózy – Vodu z vodovodu Zdarma, Potřebuji dovolenou



# č. 8 Muži bez peněz



Nic nekupují, ale nějak se chovají a víme co mají rádi

# č. 8 Muži bez peněz

insight: Hledám, co bych si koupil, ale nemám peníze

mix mužů bez peněz, hodně studenti (Státní maturity, Stáhněte si zadání) z menších měst (Brno, Ostrava)

pasivní zábava: Vyžeň nudu

sázení: Tipsport, Onlajny.cz, Chance

technické zájmy: Datart, Mobilmania, Asus, Škoda, Peugeot, Ford, Hyundai, Autosalon TV Prima

politika: Česká pirátská strana, Paroubek na Mars, Stydím se za vládu ČR

# Facebook normalized distance ...

... příběh jedné hypotézy a problémy, které jí provázejí  
a provazely

# Google distance

počítá se sémantická vzdálenost

autory jsou Rudi Cilibrasi a Paul M. B. Vitanyi

podobné věci sdílí stejné vlastnosti

tudíž se o nich mluví častěji dohromady

dvě reprezentace jsou si tím podobnější, čím méně složitých změn je třeba k převodu jedné v druhou



**NGD je vyjádřena vzorcem:**

$m = \log_{10}(\text{počet všech indexovaných stránek});$

$f_x = \log_{10}(\text{počet výsledků pro slovo } X);$

$f_y = \log_{10}(\text{počet výsledků pro slovo } Y);$

$f_{xy} = \log_{10}(\text{počet výsledků pro slovo } X \text{ a } Y);$

$GND = ((\max(f_x, f_y) - f_{xy}) / (m - \min(f_x, f_y)))$

1 apple 2 microsoft 3 bmw 4 chrysler 5 toyota

Additional keywords set (upto five, optional)

1 2 3 4 5

Restrict to domain (optional)

Domain or TLD

Example keyword set: [apple](#), [microsoft](#), [bmw](#), [chrysler](#), [toyota](#)

Result matrix

	apple	microsoft	bmw	chrysler	toyota
apple		0.11519974	0.77006889	0.85777333	0.54669557
microsoft	0.11519974		0.84528119	0.94983362	0.73886136
bmw	0.77006889	0.84528119		0.06765933	0.28139269
chrysler	0.85777333	0.94983362	0.06765933		0.42550257
toyota	0.54669557	0.73886136	0.28139269	0.42550257	

Lower numbers mean higher probability of keyword co-occurrence

 [Tweet these results](#)

<http://www.mechanicalcinderella.com/>

# FND

Facebook normalized distance vychází z Google normalized distance

Lidé komentují na stránkách kandidátů, vůči kterým se především pozitivně, ale i negativně vymezují.

Pohybují se v určitých myšlenkovinách rovinách či diskurzech

**FND je vyjádřena vzorcem:**

$m = \log_{10}(\text{počet všech českých účtů});$

$f_x = \log_{10}(\text{počet komentujících na stránce X});$

$f_y = \log_{10}(\text{počet komentujících na stránce Y});$

$f_{xy} = \log_{10}(\text{počet společných komentujících});$

$$\text{FND} = ((\max(f_x, f_y) - f_{xy}) / (m - \min(f_x, f_y)))$$

# Aplikace

bud' objevování příbuzných stránek a témat  
či k mapování vztahů mezi stránkami

# Vyhledávání

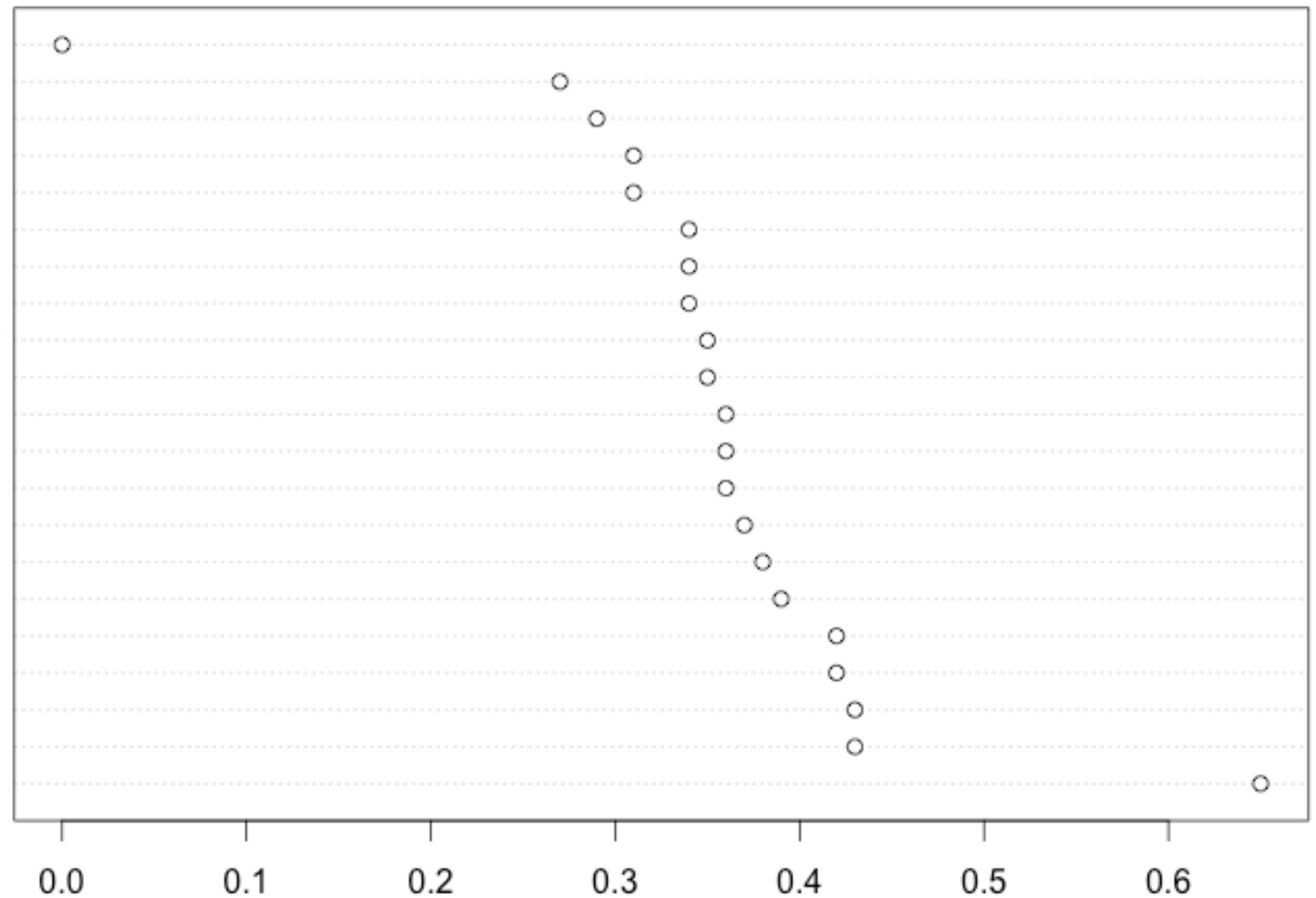
Najdi mi stránky, které mají blízko k mojí



NGD	unique count	common count	common pct	page ID	page
	3377	<a href="#">3377</a>	100.0	286483085526	<a href="#">T-Mobile CZ</a>
0.32	4923	<a href="#">519</a>	15.37	176688316811	<a href="#">Vodafone CZ</a>
0.4	8560	<a href="#">515</a>	15.25	120461304659141	<a href="#">O2 CZ</a>
0.52	4327	<a href="#">116</a>	3.44	253229733343	<a href="#">TV Nova (oficiální)</a>
0.53	4436	<a href="#">112</a>	3.32	137067469008	<a href="#">ČT24</a>
0.41	677	<a href="#">105</a>	3.11	275389913958	<a href="#">Nechceme předražené mobilní volání a služby</a>
0.43	1069	<a href="#">104</a>	3.08	115294422375	<a href="#">Mobility &amp; MobilMania.cz</a>
0.47	1454	<a href="#">90</a>	2.67	123268337716736	<a href="#">Nokia Česká republika</a>
0.54	3778	<a href="#">88</a>	2.61	142069753862	<a href="#">TN.CZ</a>
0.56	4133	<a href="#">87</a>	2.58	97718196244	<a href="#">Prima COOL</a>
0.47	1235	<a href="#">80</a>	2.37	90594467955	<a href="#">Česká spořitelna</a>
0.8	19710	<a href="#">78</a>	2.31	352840478078551	<a href="#">True story bro [CZ and SK]</a>
0.48	1257	<a href="#">78</a>	2.31	273531189144	<a href="#">Sony Ericsson Czech Republic</a>
0.54	3063	<a href="#">78</a>	2.31	51828152685	<a href="#">Česká televize</a>
0.47	1027	<a href="#">76</a>	2.25	188676337835075	<a href="#">Air Bank</a>
0.8	19219	<a href="#">74</a>	2.19	109204525794230	<a href="#">"Au moje koule!"</a>
0.7	8513	<a href="#">65</a>	1.92	115991108439011	<a href="#">EVROPA 2</a>
0.72	9153	<a href="#">63</a>	1.87	111003257404	<a href="#">Fajn Radio</a>
0.52	1662	<a href="#">63</a>	1.87	191987746428	<a href="#">Ostře sledované žluté vlaky</a>
0.56	2818	<a href="#">62</a>	1.84	312917027851	<a href="#">Slevomat</a>
0.48	712	<a href="#">56</a>	1.66	361281313901258	<a href="#">Komerční banka</a>

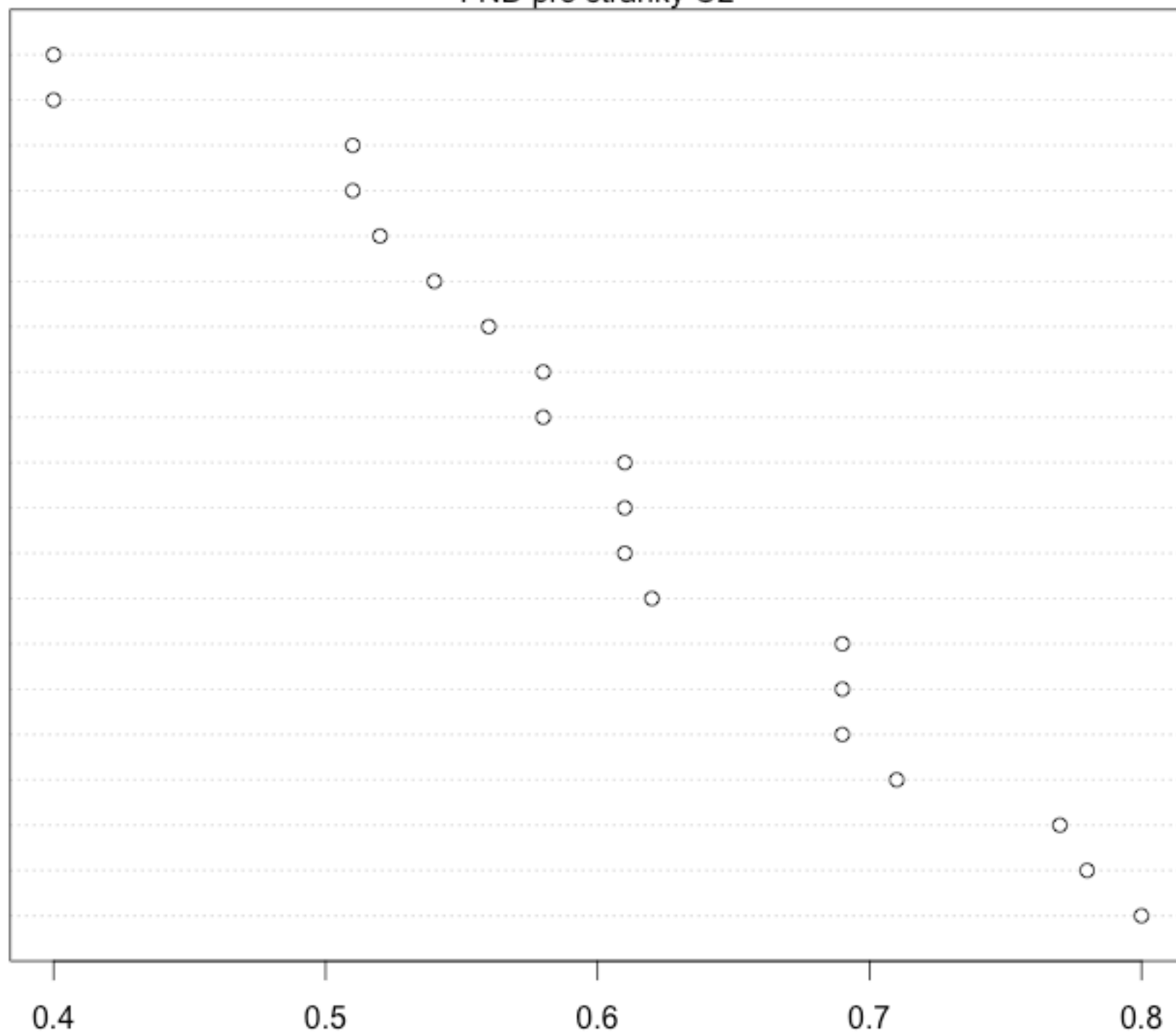
## Blízkost stránek ke stránce NIVEA Friends

NIVEA Friends  
Garnier CZ  
Maybelline CZ/SK  
Fidorka  
Milka  
Garnier SK  
Bourjois Paris v ČR  
L'Oréal Paris CZ  
Schwarzkopf Professional - Vášeň pro vlasy  
MasterCard Česká republika  
Kinder Bueno CZ  
Sedita - Chuť, která Tě potěší  
Tescoma  
Manner CZ  
BIODERMA  
Hamánek  
Avon Česká republika  
Pro dětičky  
WeBrouček  
Kafe.cz - pohoda pro ženy  
O2 CZ



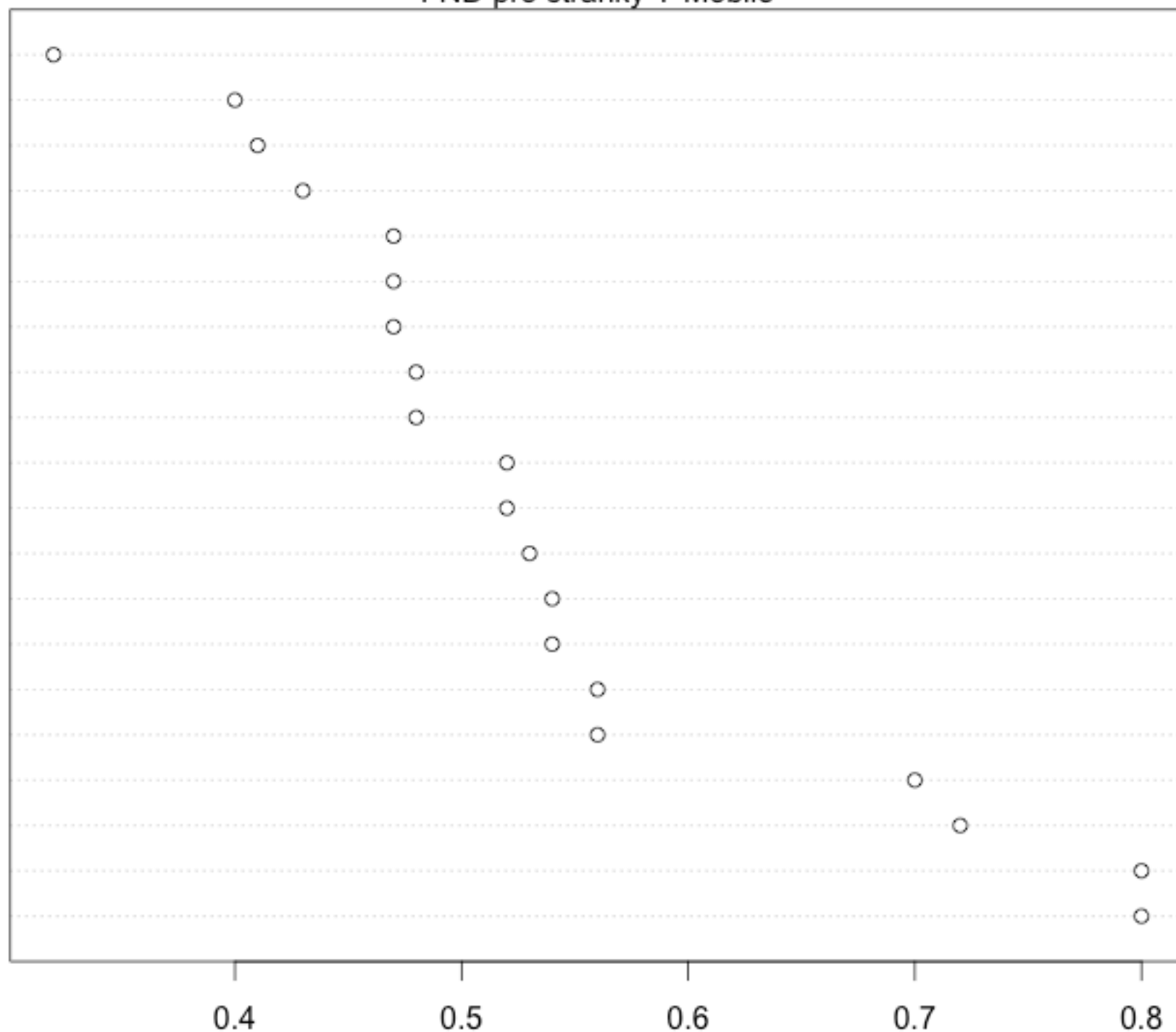
FND pro stránky O2

- T-Mobile CZ
- Vodafone CZ
- Nechceme předražené mobilní volání a služby
- Mobility & MobilMania.cz
- Nokia Česká republika
- Česká spořitelna
- Sony Ericsson Czech Republic
- Ostře sledované žluté vlaky
- TV Nova (oficiální)
- Česká televize
- TN.CZ
- ČT24
- Prima COOL
- RE-PLAY
- EVROPA 2
- Fajn Radio
- Kiss Jižní Čechy
- Au moje koule!
- Rozzlobení muži
- True story bro [CZ and SK]



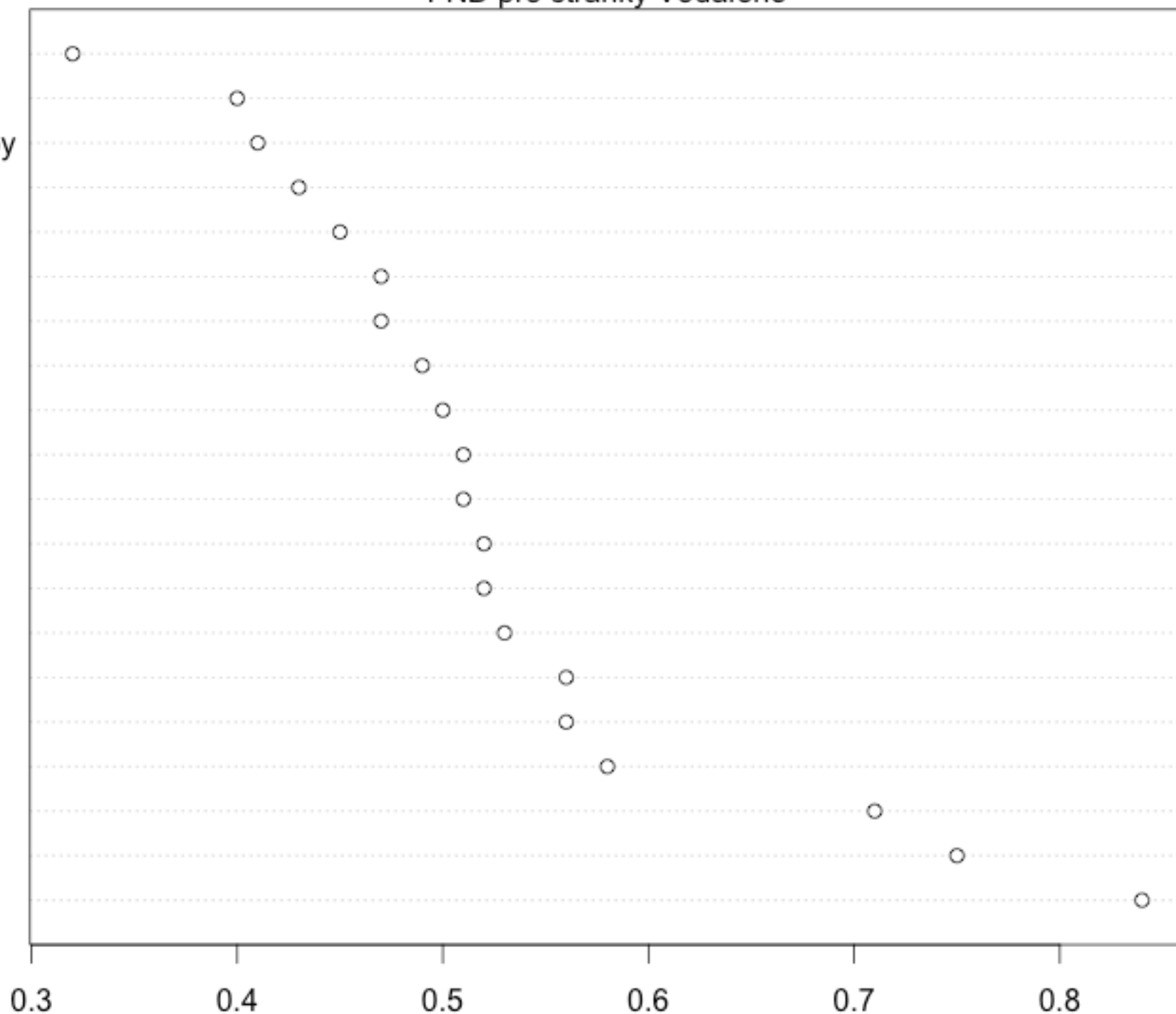
FND pro stránky T-Mobile

- Vodafone CZ
- O2 CZ
- Nechceme předražené mobilní volání a služby
- Mobility & MobilMania.cz
- Air Bank
- Česká spořitelna
- Nokia Česká republika
- Komerční banka
- Sony Ericsson Czech Republic
- Ostře sledované žluté vlaky
- TV Nova (oficiální)
- ČT24
- Česká televize
- TN.CZ
- Slevomat
- Prima COOL
- EVROPA 2
- Fajn Radio
- Au moje koule!
- True story bro [CZ and SK]

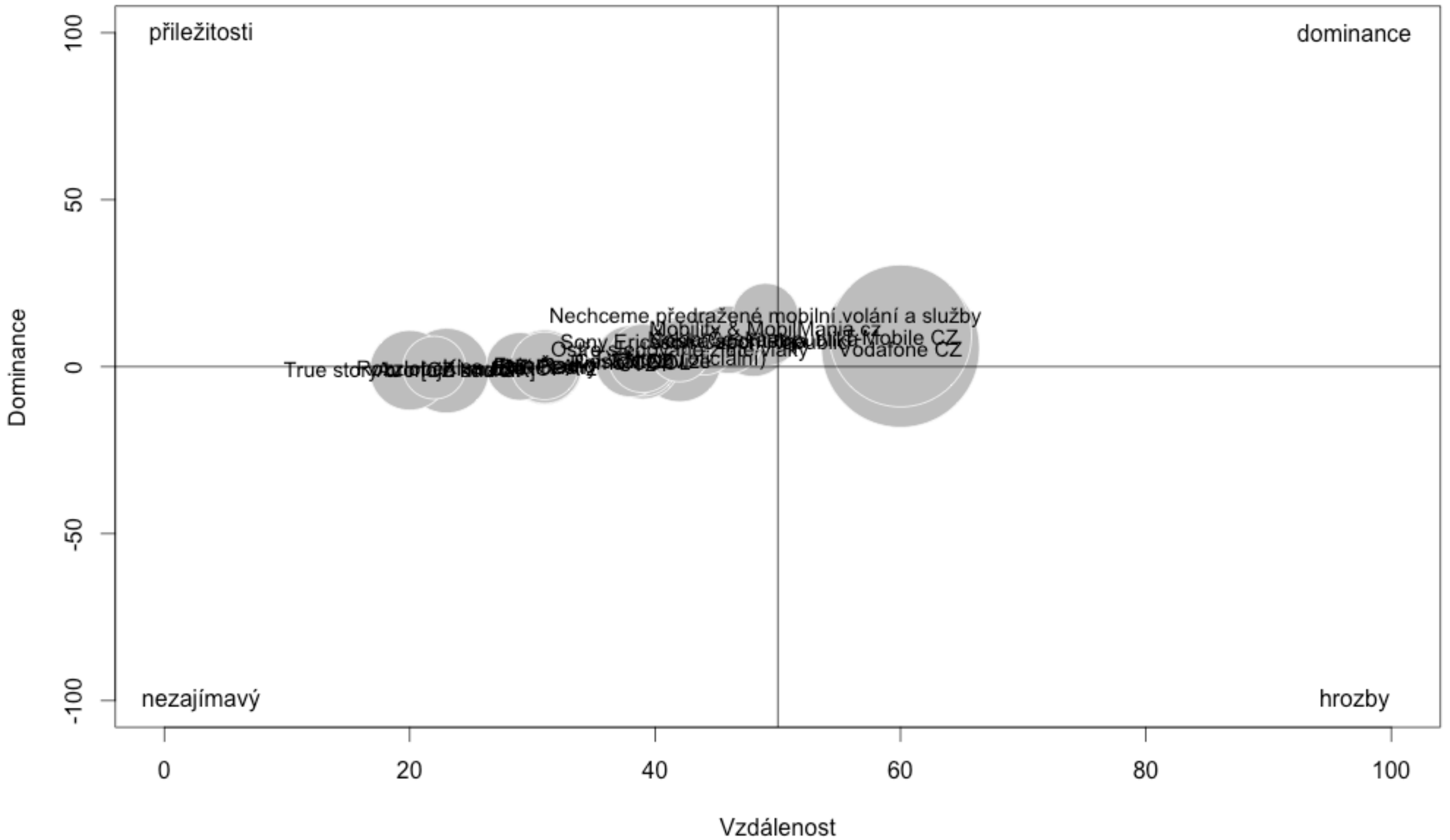


FND pro stránky Vodafone

- T-Mobile CZ
- O2 CZ
- Nechceme předražené mobilní volání a služby
- Mobility & MobilMania.cz
- Air Bank
- Sony Ericsson Czech Republic
- Nokia Česká republika
- Česká spořitelna
- LEO Express
- HTC Fans Cz/Sk
- Ostře sledované žluté vlaky
- Česká televize
- ČT24
- TV Nova (oficiální)
- Slevomat
- TN.CZ
- Prima COOL
- EVROPA 2
- Au moje koule!
- True story bro [CZ and SK]

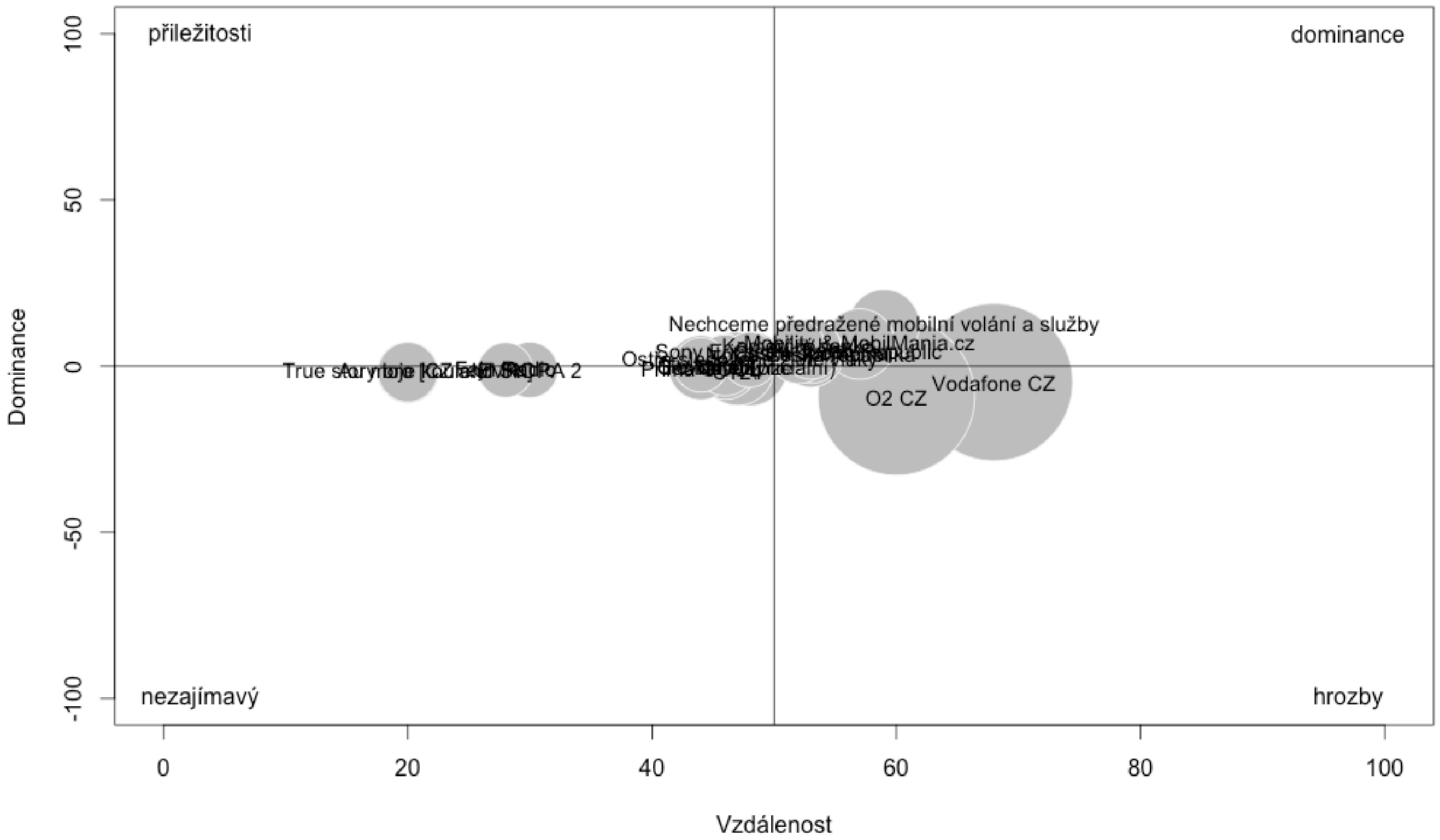


# Magic Quadrant for O2 page on Facebook





# Magic Quadrant for T-Mobile page on Facebook

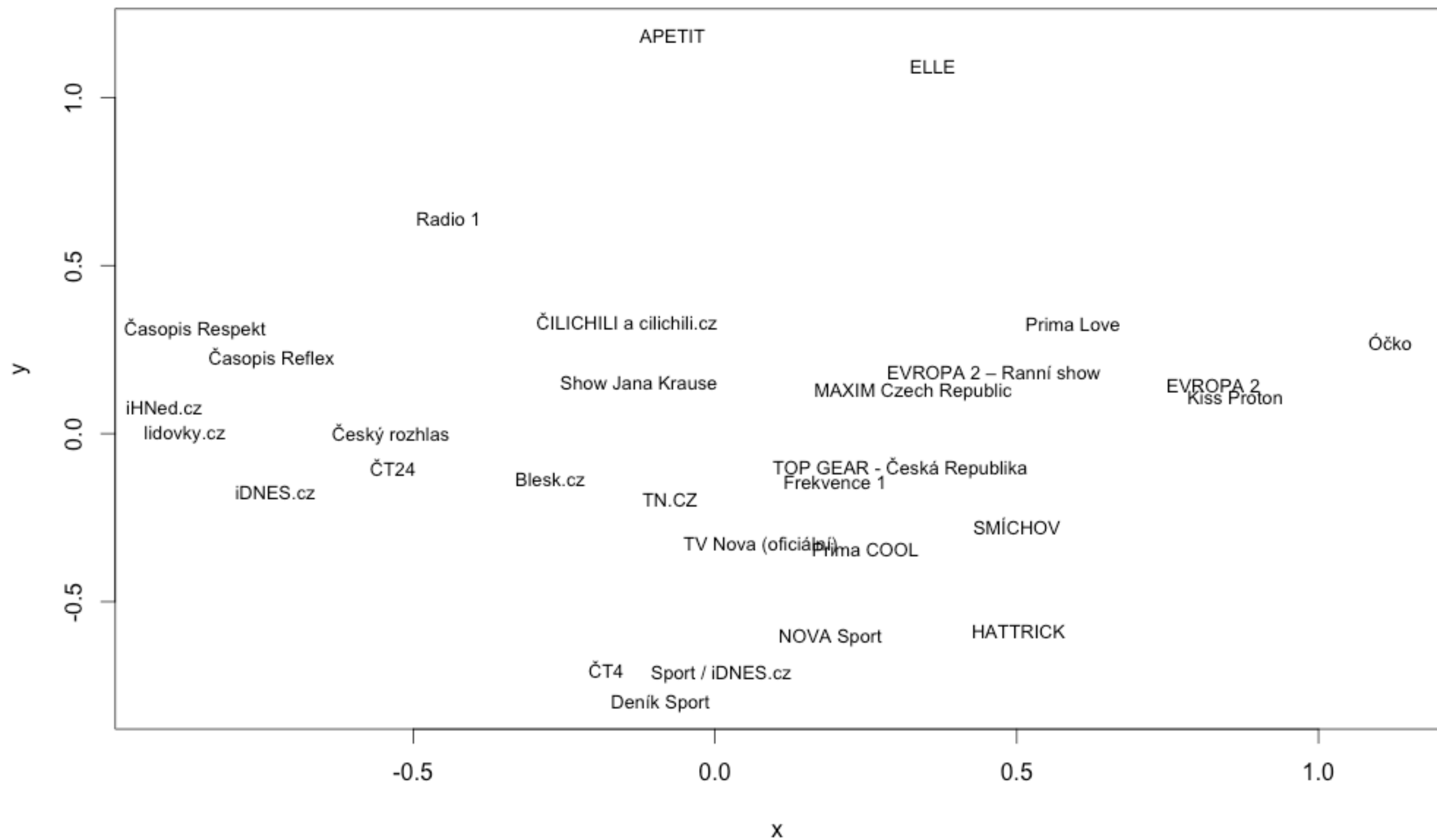


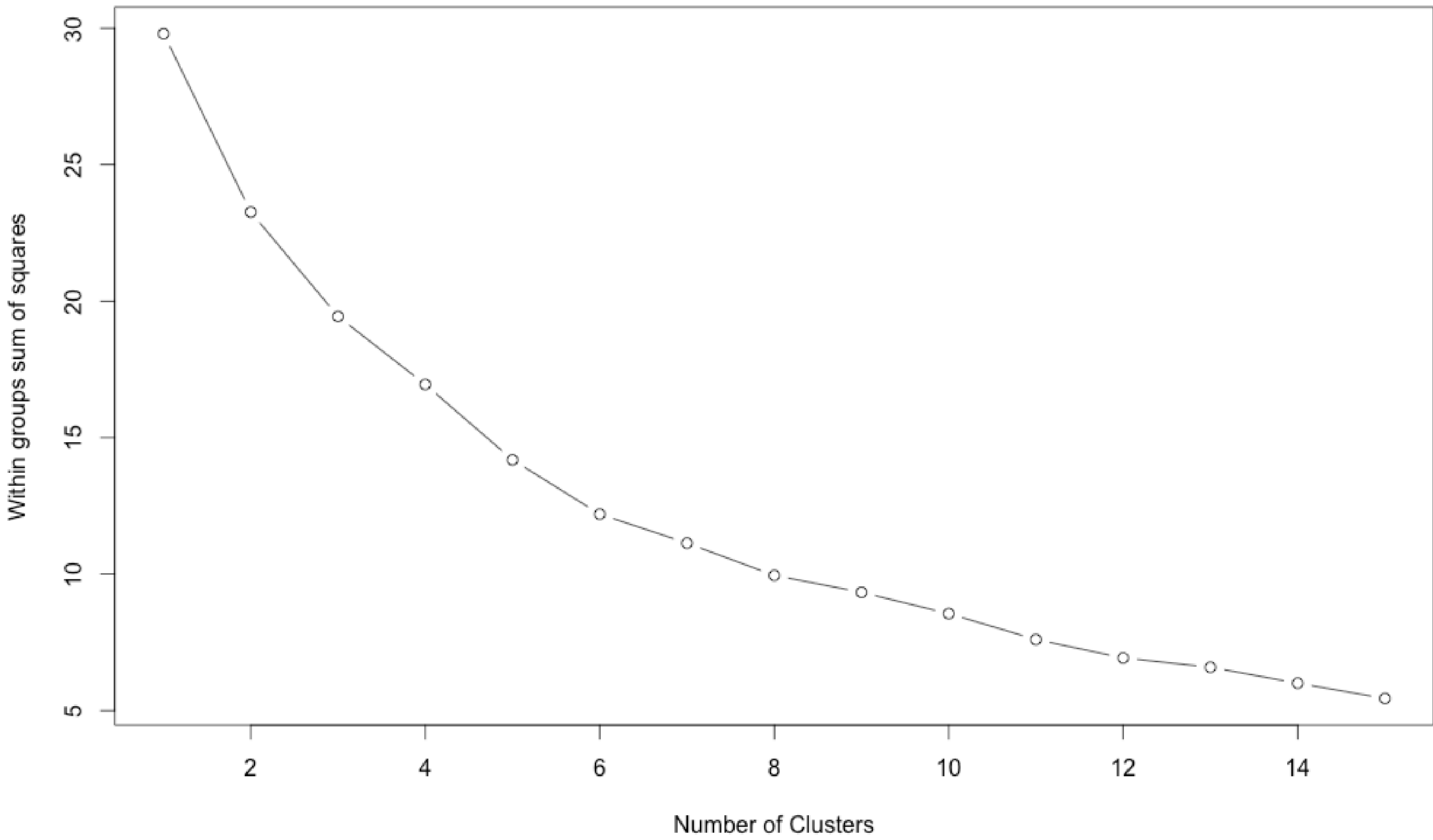
# Vztahování

Kde se nachází moje stránka mezi ostatními

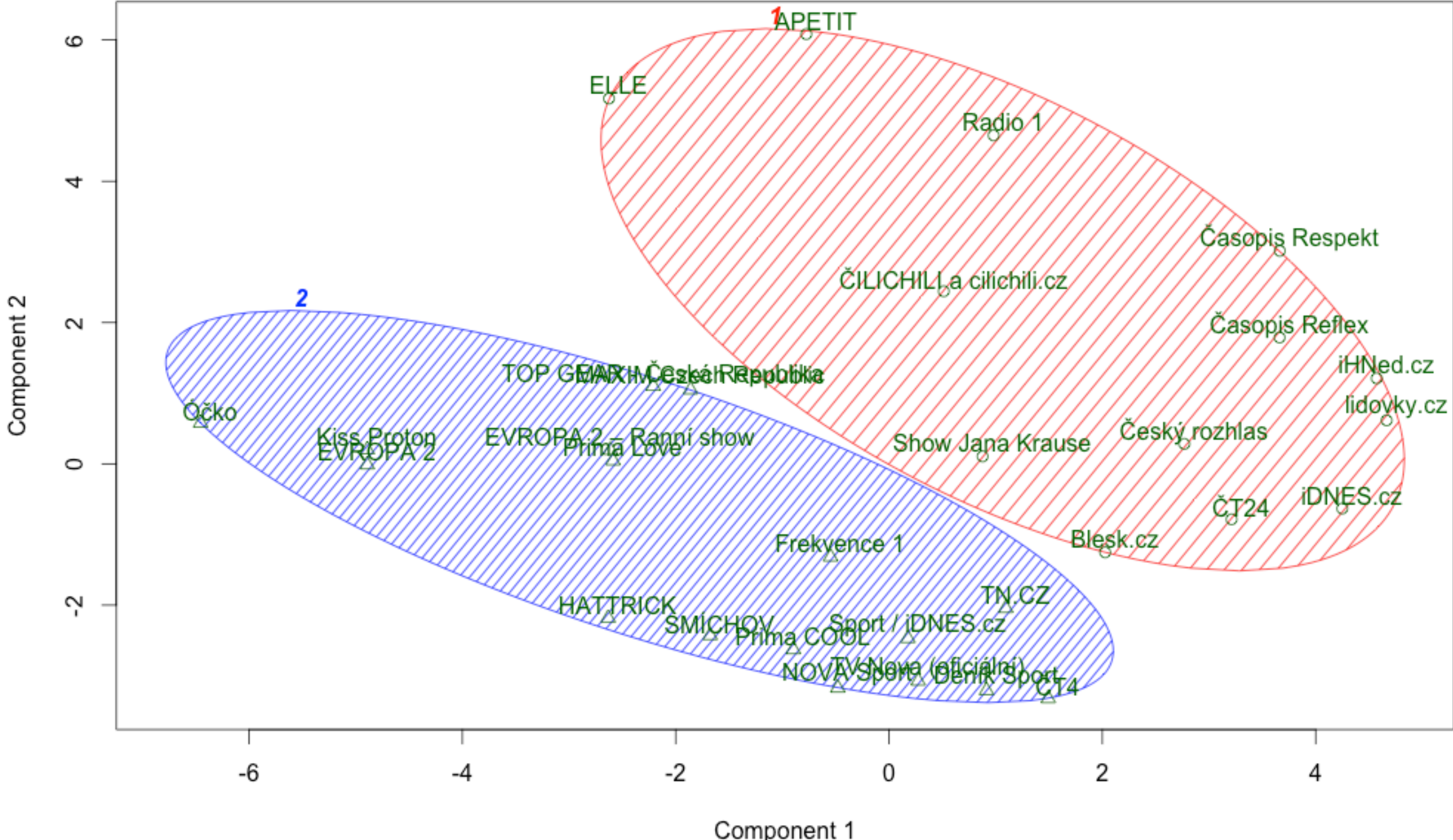
# Mediální mapa FB

Jaké typy stránek máme mezi FB stránkami serveru





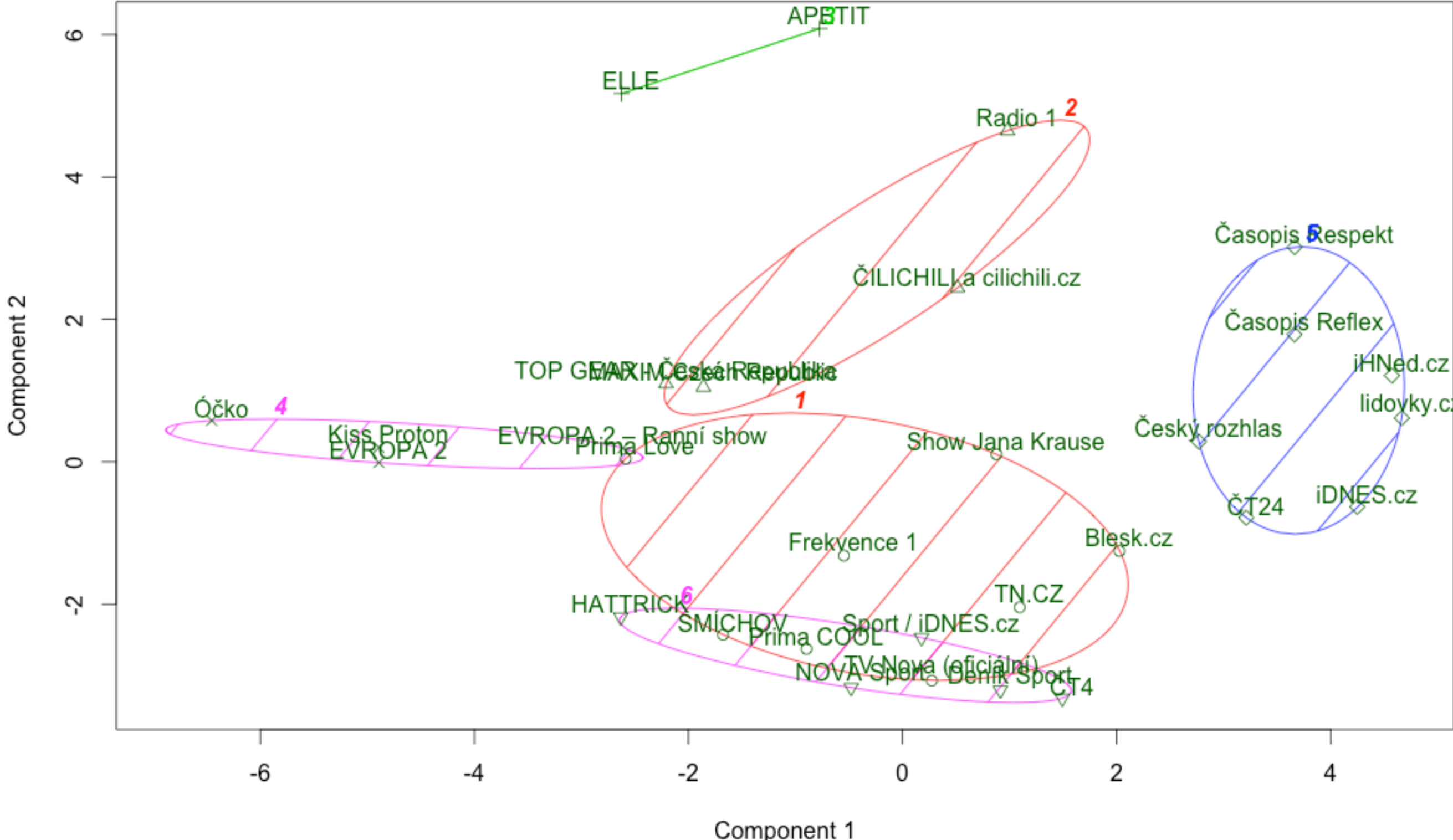
# CLUSPLOT( tbl )



These two components explain 47.47 % of the point variability.

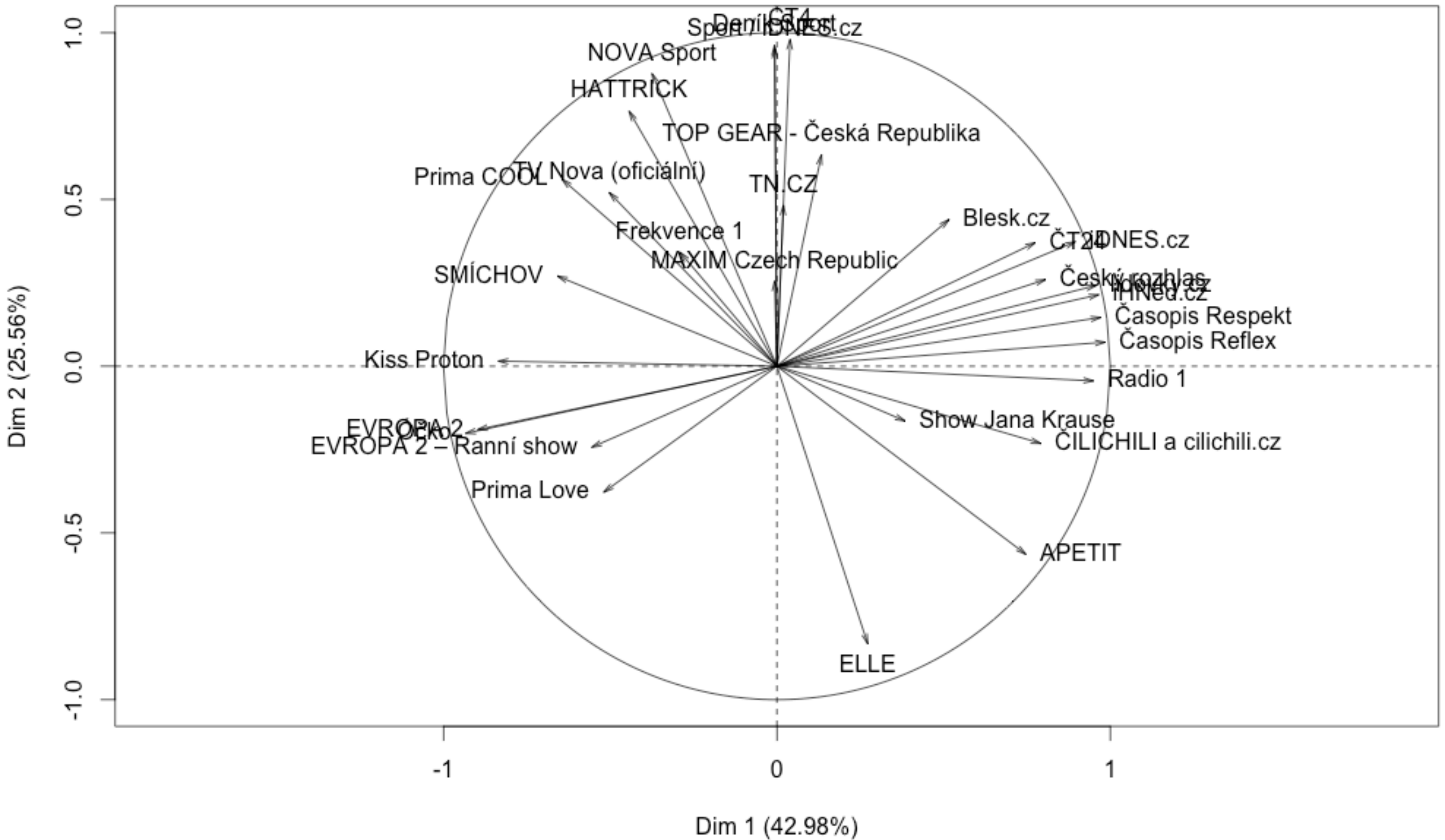


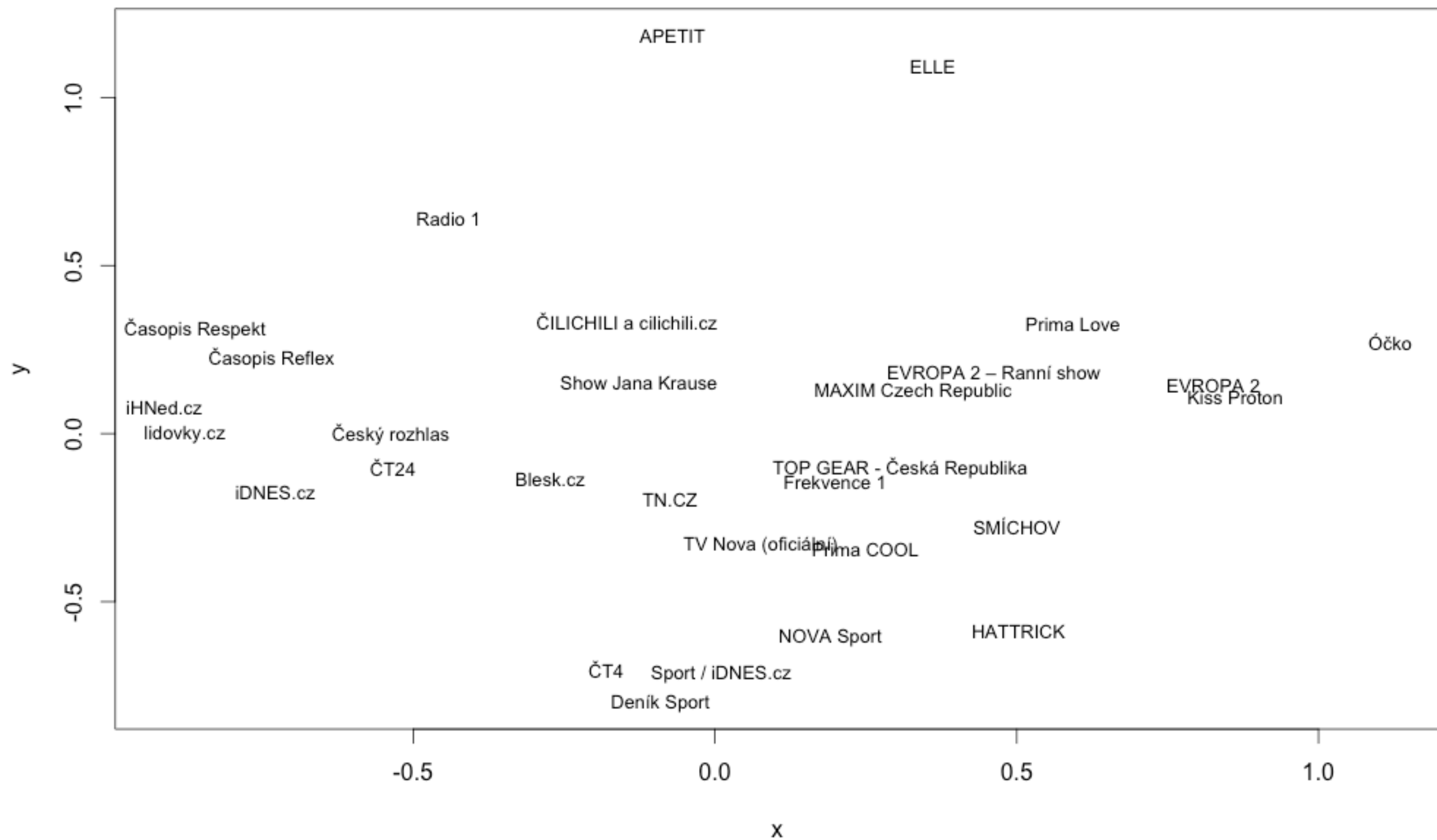
# CLUSPLOT( tbl )



These two components explain 47.47 % of the point variability.

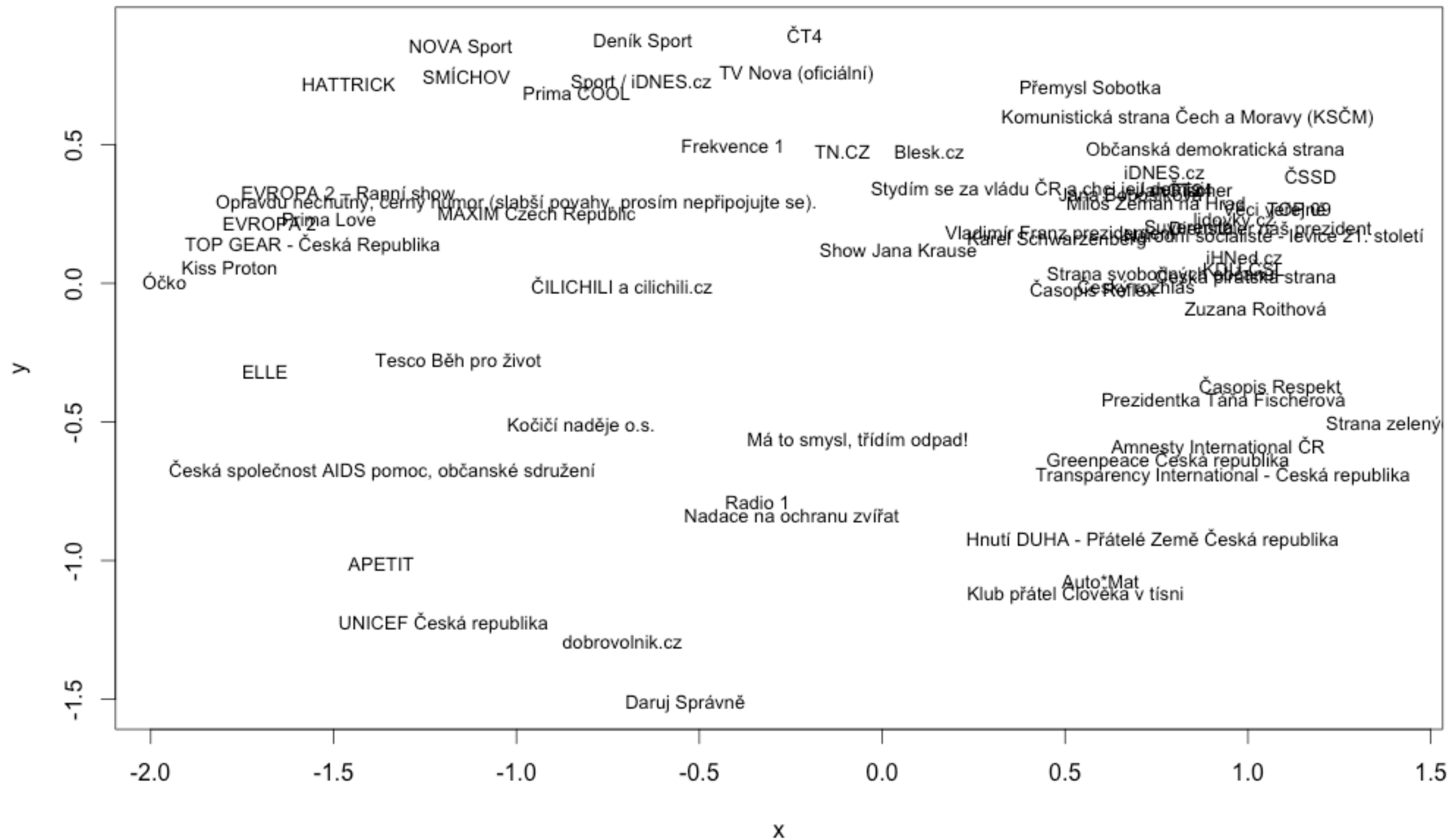
Variables factor map (PCA)

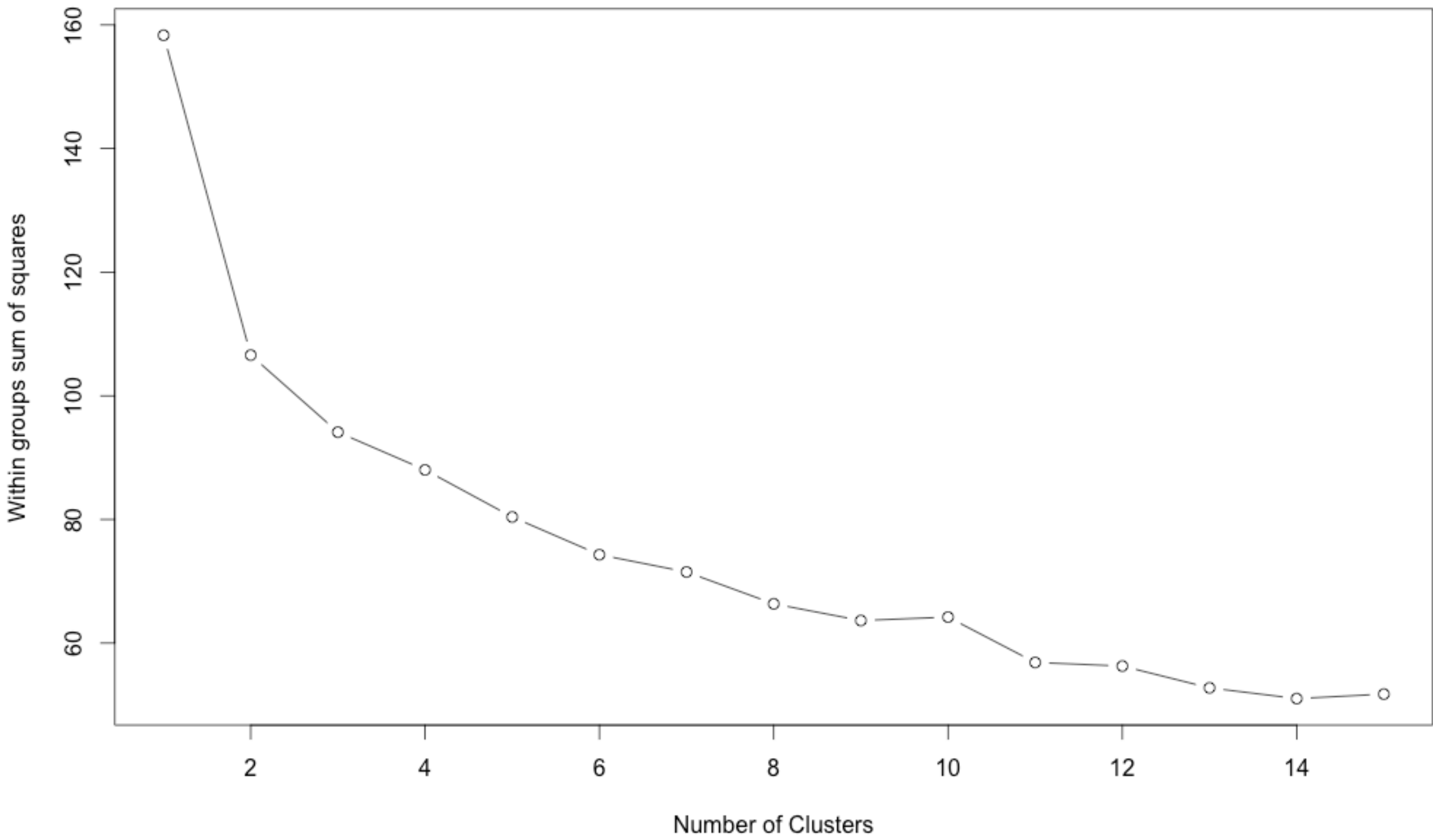




# Velká mapa stránek

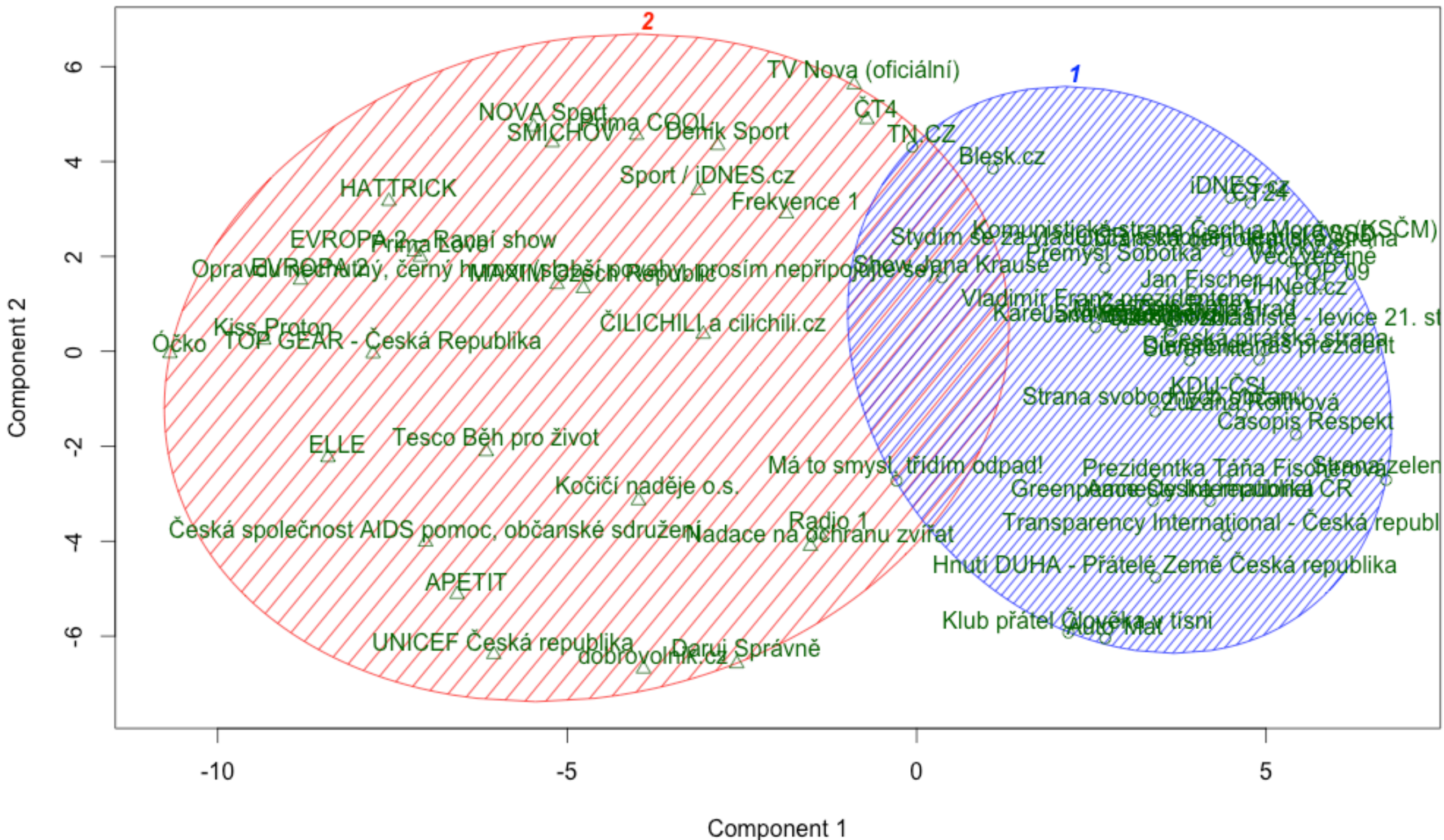
Neziskovky, strany, média....







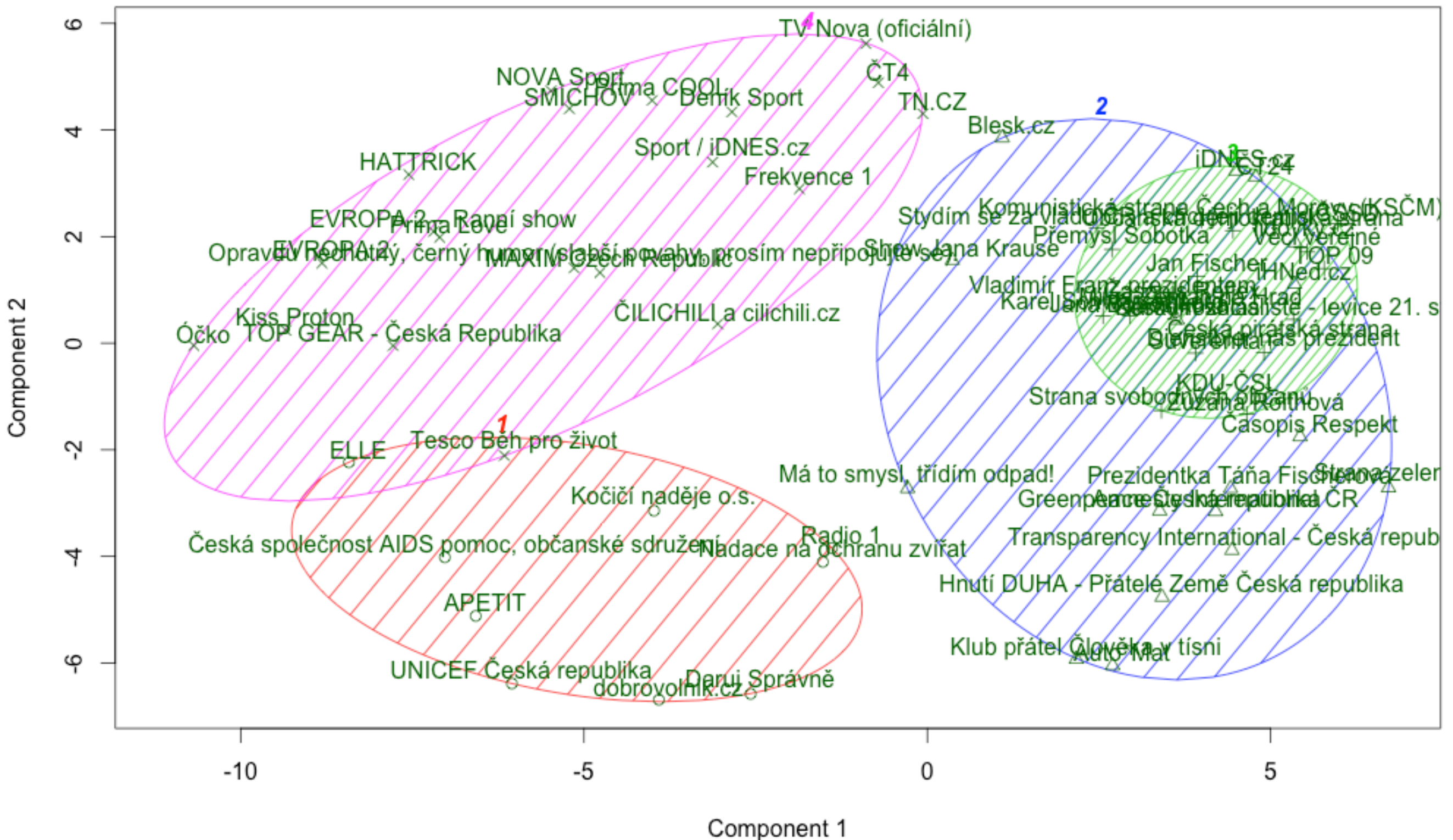
# CLUSPLOT( tbl )



These two components explain 51.1 % of the point variability.

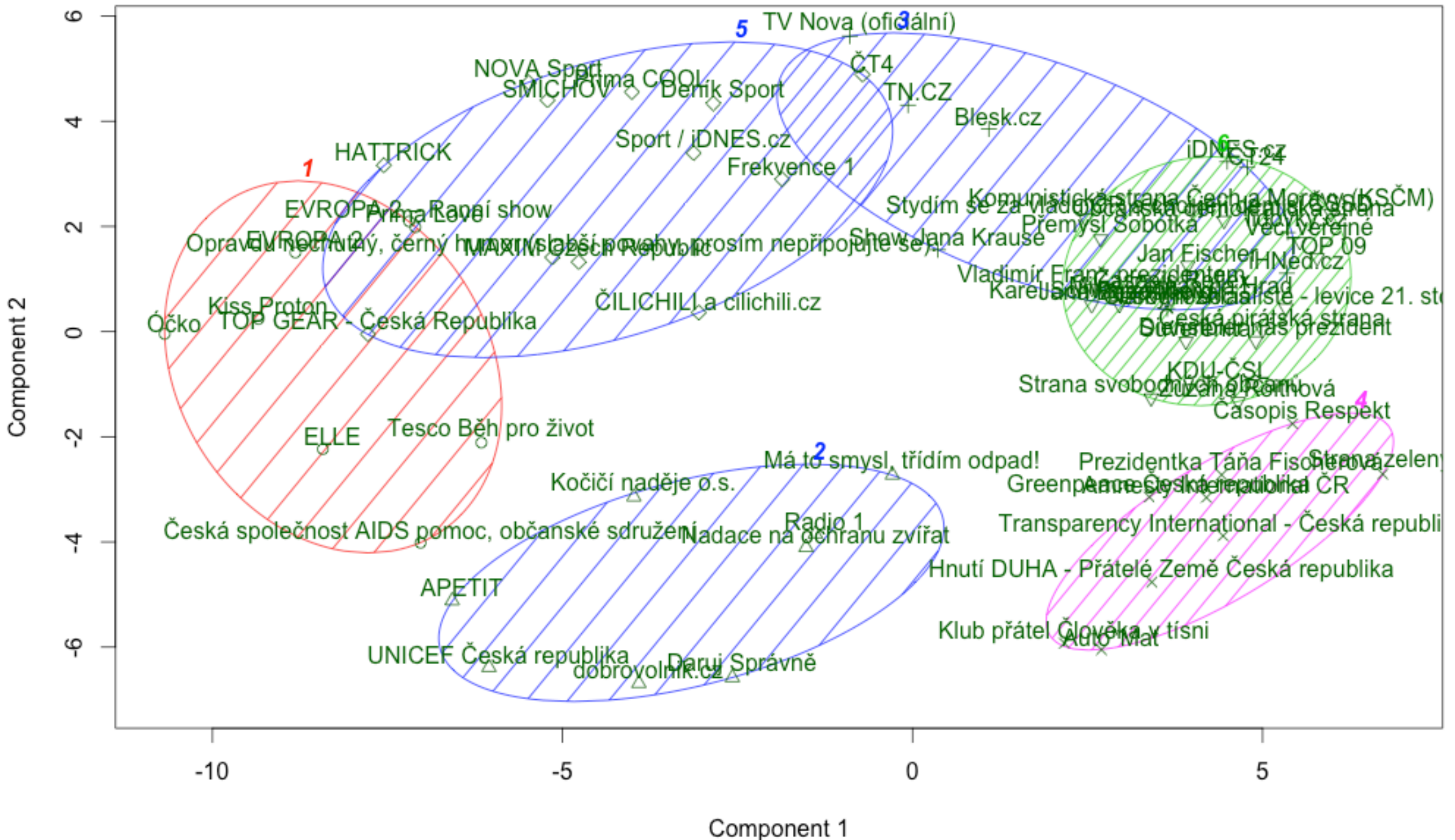


# CLUSPLOT( tbl )



These two components explain 51.1 % of the point variability.

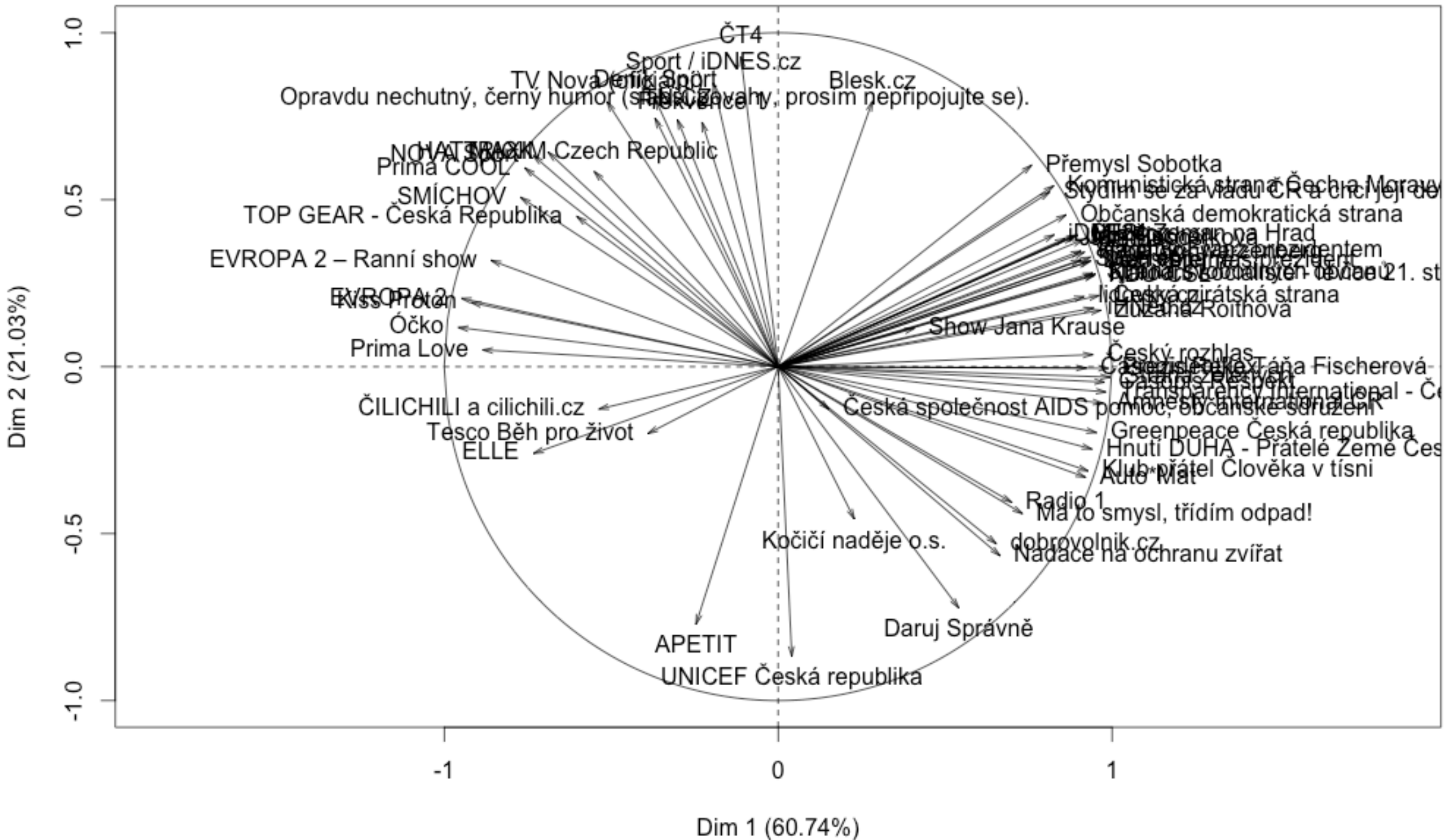
# CLUSPLOT( tbl )

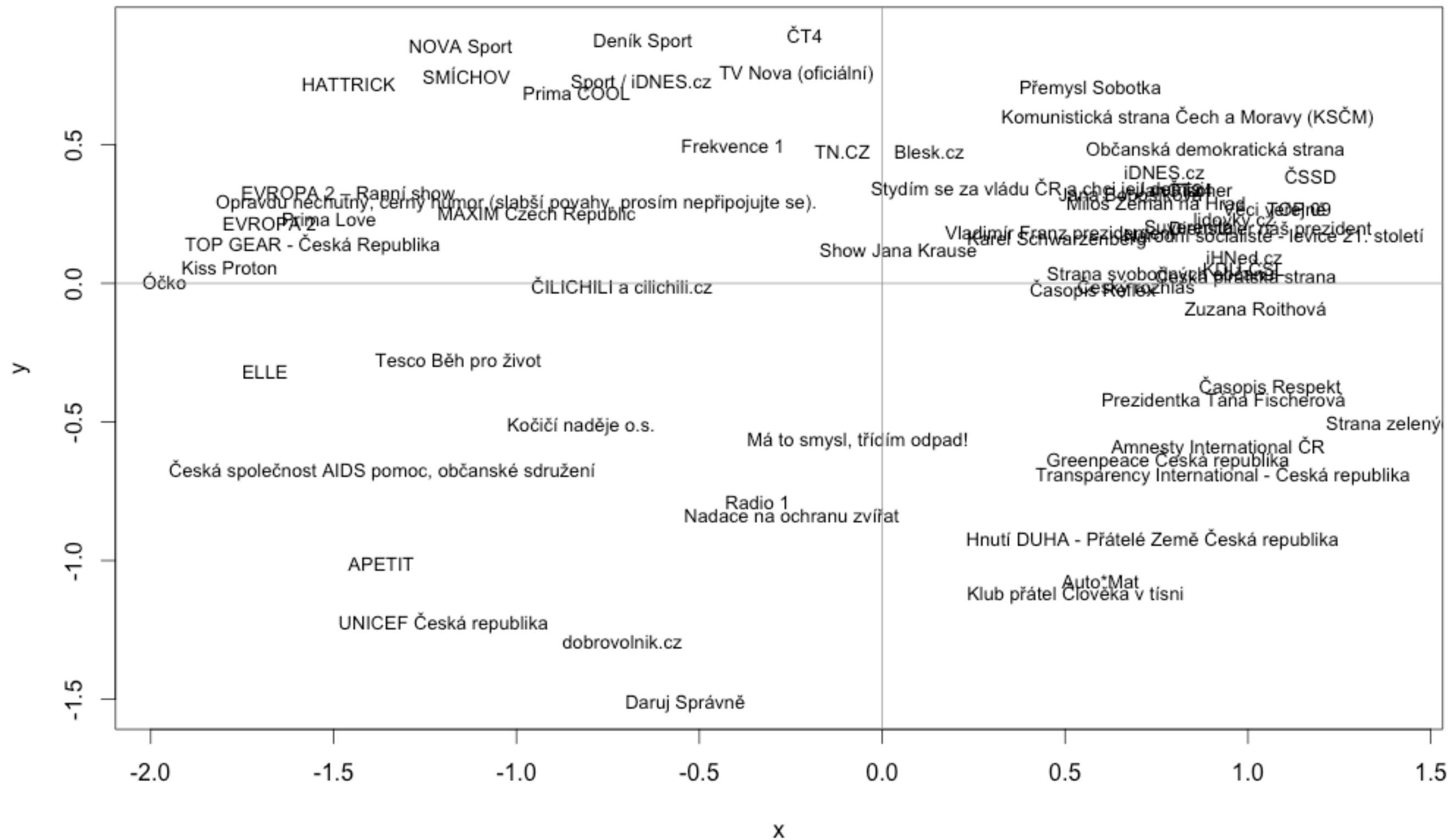


These two components explain 51.1 % of the point variability.



Variables factor map (PCA)



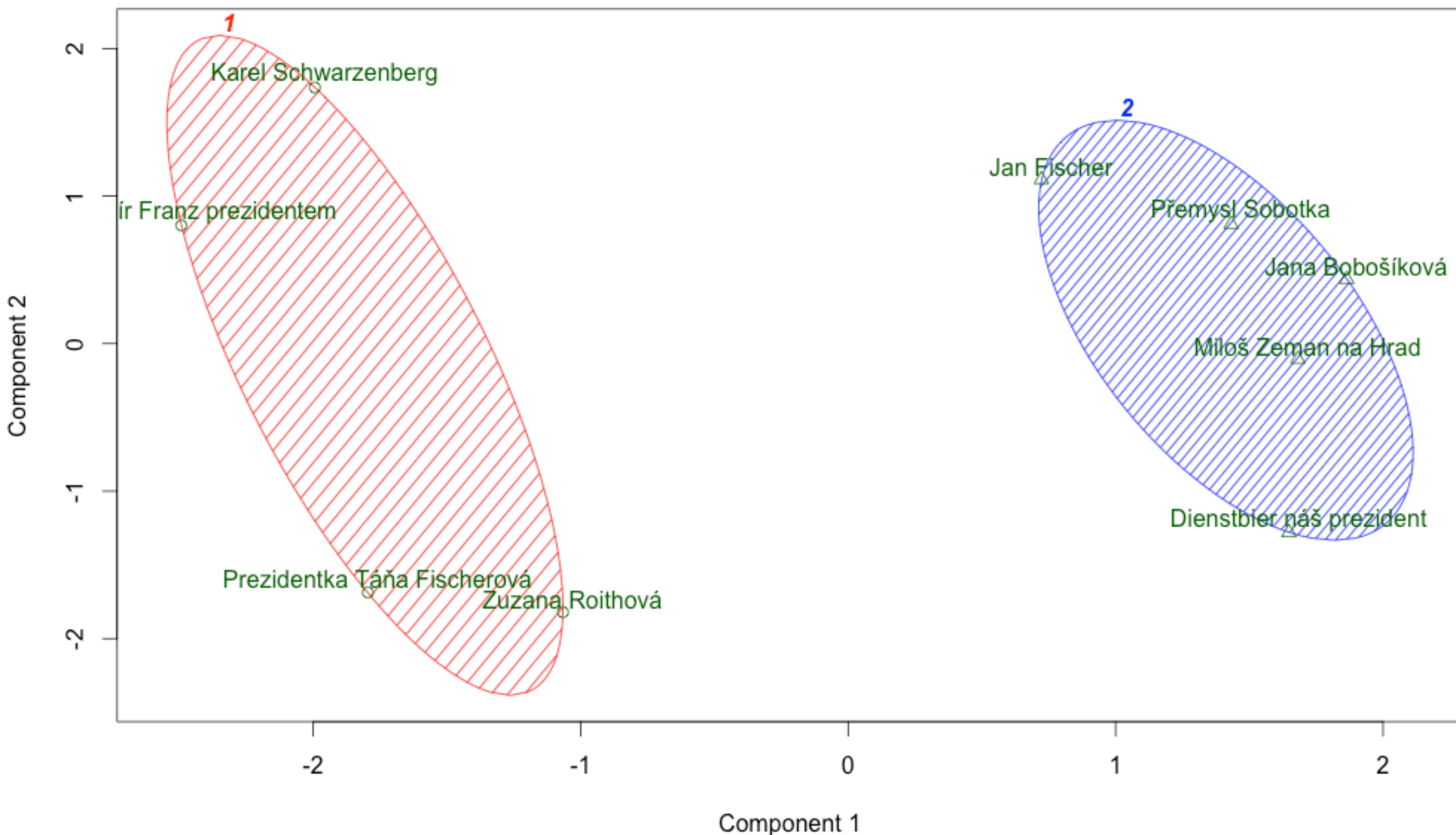


# Případová studie

Prezidentské volby 2013

# Mapa prezidentských kandidátů 14. 1. 2013

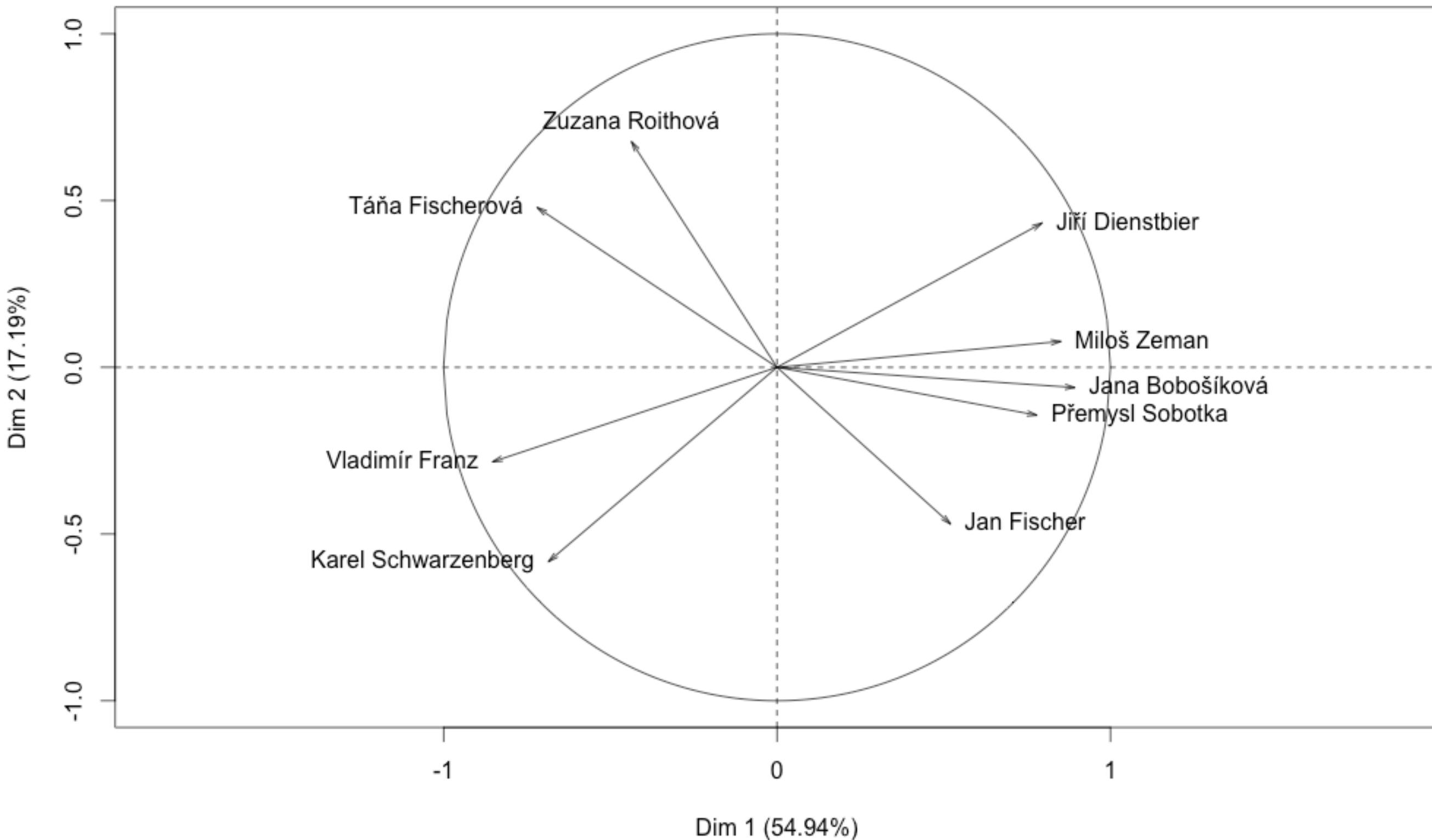
CLUSPLOT( tbl )



These two components explain 49.01 % of the point variability.

# Mapa prezidentských kandidátů 14. 1. 2013

Variables factor map (PCA)



# Hlavní úskalí

velikosti korpus a průniku (vzorec funguje “vždy”)

validace

úplnost korpusu při objevování

relativní drahost

vývoj v čase



# Co nás čeká

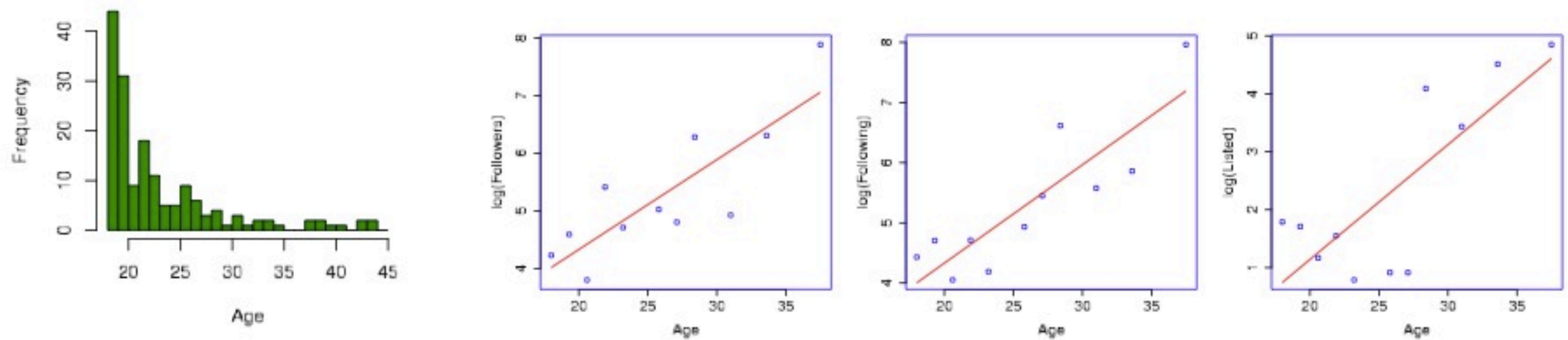
aneb na čem teď pracujeme

**Tabuľka 1.3***Prehľad charakteristík obsiahnutých v päťfaktorovom modeli osobnosti*

Opis viazaný na vysoké skóre	Konvenčný názov a definícia faktoru	Opis viazaný na nízke skóre
spoločenský, aktívny, zhovorčivý, družný, priebojný, smelý	<b>Extraverzia (E; Extraversion)</b> Zisťuje kvantitu interpersonálnych interakcií, úroveň aktivácia a stimulácie.	uzavretý, plachý, málovravný, nesmelý, samotársky, tichý
dobrosrdečný, láskavý, znášateľný, srdečný, skromný, zmierlivý	<b>Prívetivosť (A; Agreeableness)</b> Zisťuje kvalitu interpersonálnej orientácie na kontinuu od súcitenia po nepriateľskosť v myšlienkach, pocitoch i činoch.	panovačný, útočný, pomstychtivý, bezcitný, despotický, konfliktný
spoľahlivý, pracovitý, presný,	<b>Svedomitosť (C; Conscientiousness)</b> Zisťuje mieru motivácie a vytrvalosti na cieľ zameraného správania. Odlišuje spoľahlivých, na	bezcieľný, nedbalý, lenivý, nesvedomitý,

<p>poriadkumilovný, zodpovedný, starostlivý</p>	<p>seba náročných ľudí od tých, ktorí sú ľahostajní a nedbalí.</p>	<p>chaotický, nevytrvalý</p>
<p>napätý, nepokojný, labilný, neistý, vznetlivý, popudlivý</p>	<p><b>Emocionálna (ne)stabilita / Neuroticizmus (N; Neuroticism)</b>  Odlišuje jedincov náchylných k psychickému vyčerpaniu a ťažko zvládajúcich psychickú záťaž od jedincov vyrovnaných a odolných voči psychickému vyčerpaniu.</p>	<p>pokojný, uvoľnený, vyrovnaný, stabilný, sebaistý, nezdolný</p>
<p>zvedavý, originálny, tvorivý, obrazotvorný, inteligentný, kultivovaný</p>	<p><b>Otvorenosť voči skúsenosti / Intelekt / Kultúra (I; Intellect)</b>  Zisťuje aktívne vyhľadávanie nových skúseností, zážitkov, toleranciu neznámeho.</p>	<p>konvenčný, pragmatický, realistický, neprispôsobivý, neinteligentný, nevzdelaný</p>





(a) Distributions of Age.

(b) User activity varies with age.

Fig. 3. Age distribution and effects on Twitter activity

Trait	Listeners <i>log(Following)</i>	Popular <i>log(Followers)</i>	Highly-read <i>log(Listed)</i>	Influential <i>Klout</i>	Influential <i>log(TIME)</i>
O	0.05	0.05	<b>0.17*</b>	0.13	0.00
C	0.08	0.10	0.02	0.01	<b>0.18***</b>
E	<b>0.13*</b>	<b>0.15**</b>	0.09	<b>0.15*</b>	<b>0.25***</b>
A	0.07	0.02	0.03	-0.17	0.06
N	<b>-0.17**</b>	<b>-0.19***</b>	-0.03	<b>-0.03*</b>	<b>-0.20***</b>
<i>log(Age)</i>	<b>0.28*</b>	<b>0.37*</b>	0.13	0.05	<b>0.39*</b>
Male	-0.05	-0.05	-0.05	-0.04	0.01

**Děkuju za pozornost**

@josefslerka