

5 May 2021

Las Médulas. Spanish Roman over ground gold mine

Authors:

Ramón Alvarez-Esteban (1) (Maintainer)

Mónica Bécue-Bertaut (2)

Josep-Anton Sánchez-Espigares (2)

Belchin Kostov (2)

(1) University of Leon / Spain. ramon.alvarez@unileon.es

(2) UPC Universitat Politècnica de Catalunya / Spain

XplorText

Multivariate statistical methods for analyze textual data

<http://www.xplorText.org/>

<https://xplorText.unileon.es/>

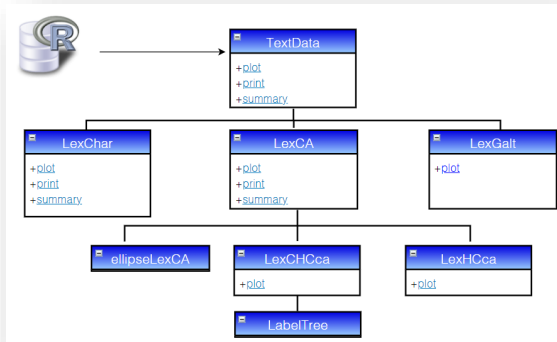


This script



R User Group

Home Flowchart Analyse textuelle avec R Textual Data Science With R Spanish political speeches



- print
- summary
- plot

Packages:

tm(>= 0.7-8)
slam(>= 0.1-48)

ggplot2
ggdendro(>= 0.1.22)
ggforce(>= 0.3.2)
ggrepel(>= 0.9.0),
graphics
gridExtra(>= 2.3),

FactoMineR, MASS, methods, stats,
utils, flexclust, flashClust

stringi(>= 1.5.3), stringr(>= 1.4.0)



Open question:

what is the most important thing in life?

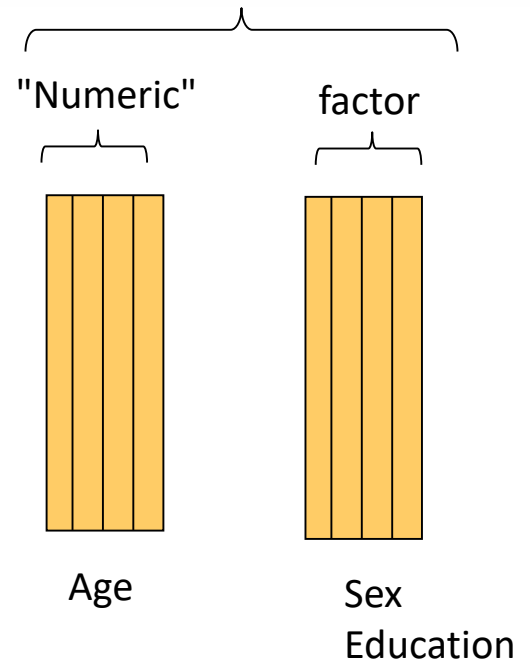
Enriching the processing of texts with contextual variables

UTF-8

Textual variable (Corpus)

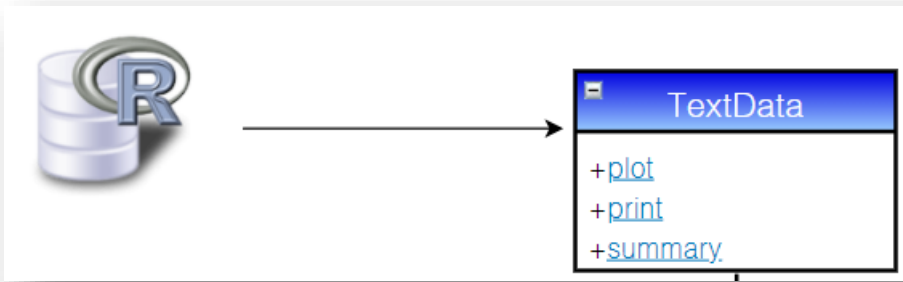
Documents

1	Good health. Happiness
2	Healthy, have enough to eat,...
.	
.	
.	
1043	job, good life, money, health



Working with dataframes (rownames for the documents)
No tibbles

UTF-8



TextData function

var.agg

context.quali

context.quantil

selDoc

tm options

lower=TRUE, remov.number,

lminword, Fmin, Dmin, Fmax

stop.word.tm=FALSE, idiom="en"

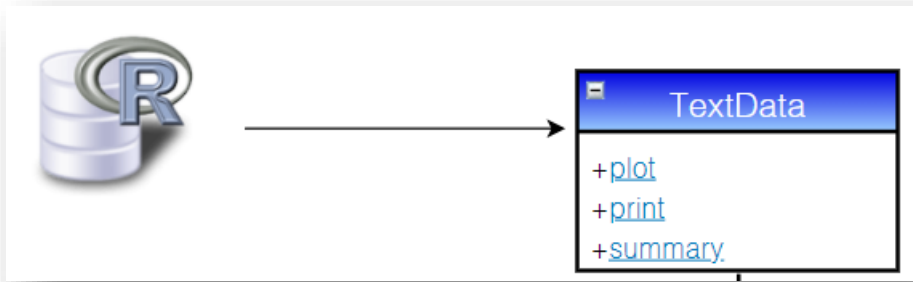
stop.word.user

segments

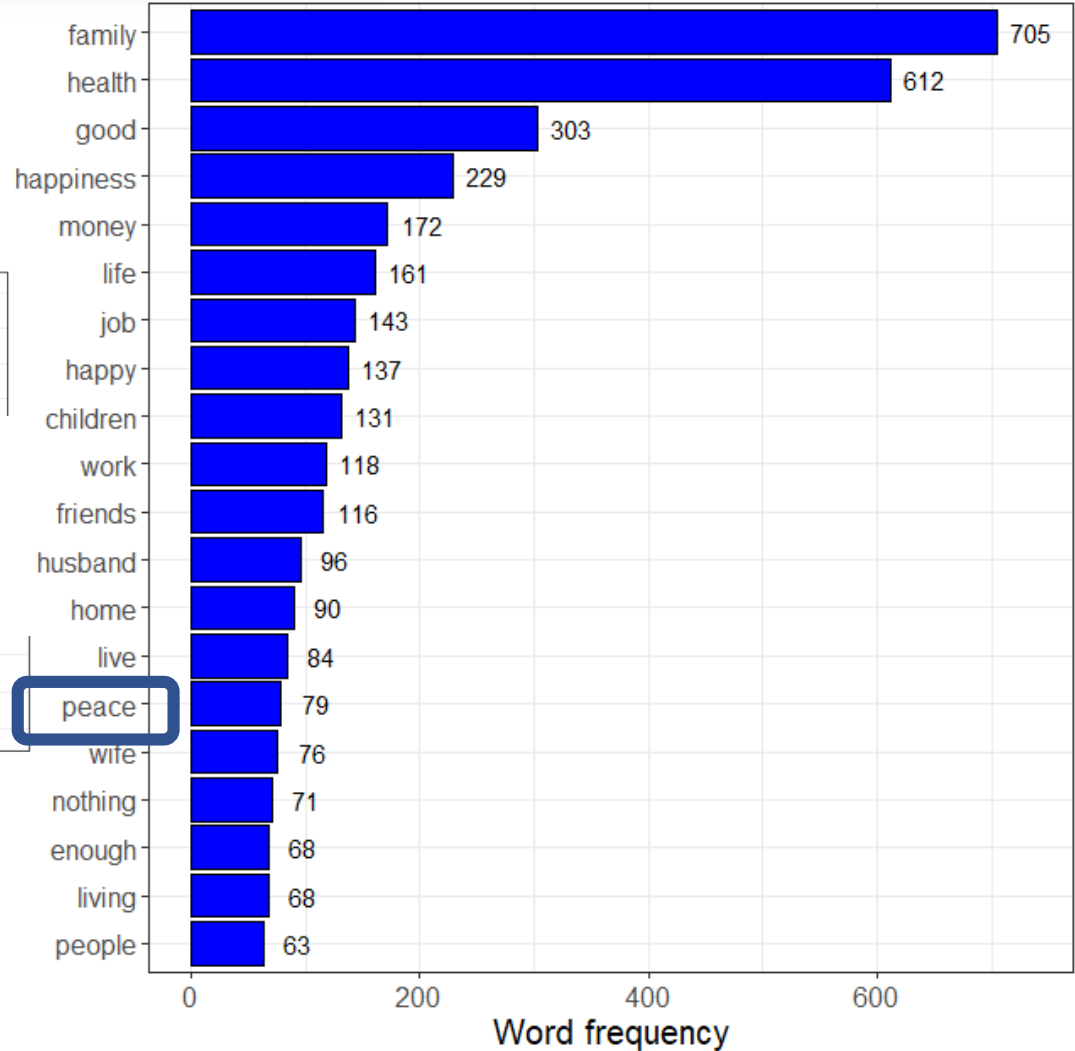


```
res.TD <-TextData(data, var.text=c(9,10), var.agg="Age12Categ", remov.number=TRUE,
  Fmin=10, Dmin=10, stop.word.tm=TRUE)
summary(res.TD)
```

	Before	After
Documents	1043.00	12.00
Occurrences	13917.00	5911.00
Words	1334.00	132.00
Mean-length	13.34	492.58
NonEmpty.Docs	1040.00	12.00
NonEmpty.Mean-length	13.38	492.58

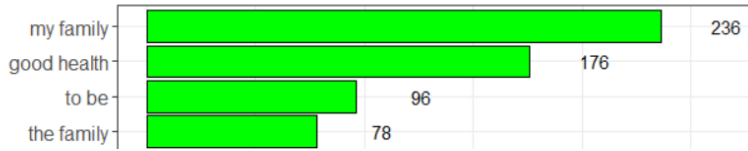


20 most frequent words

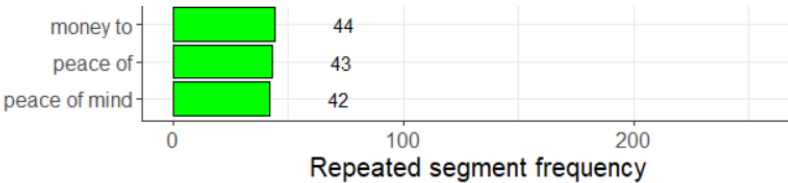


```
plot(res.TD, col.fill="blue")
```

20 most frequent segments



peace of mind: 42

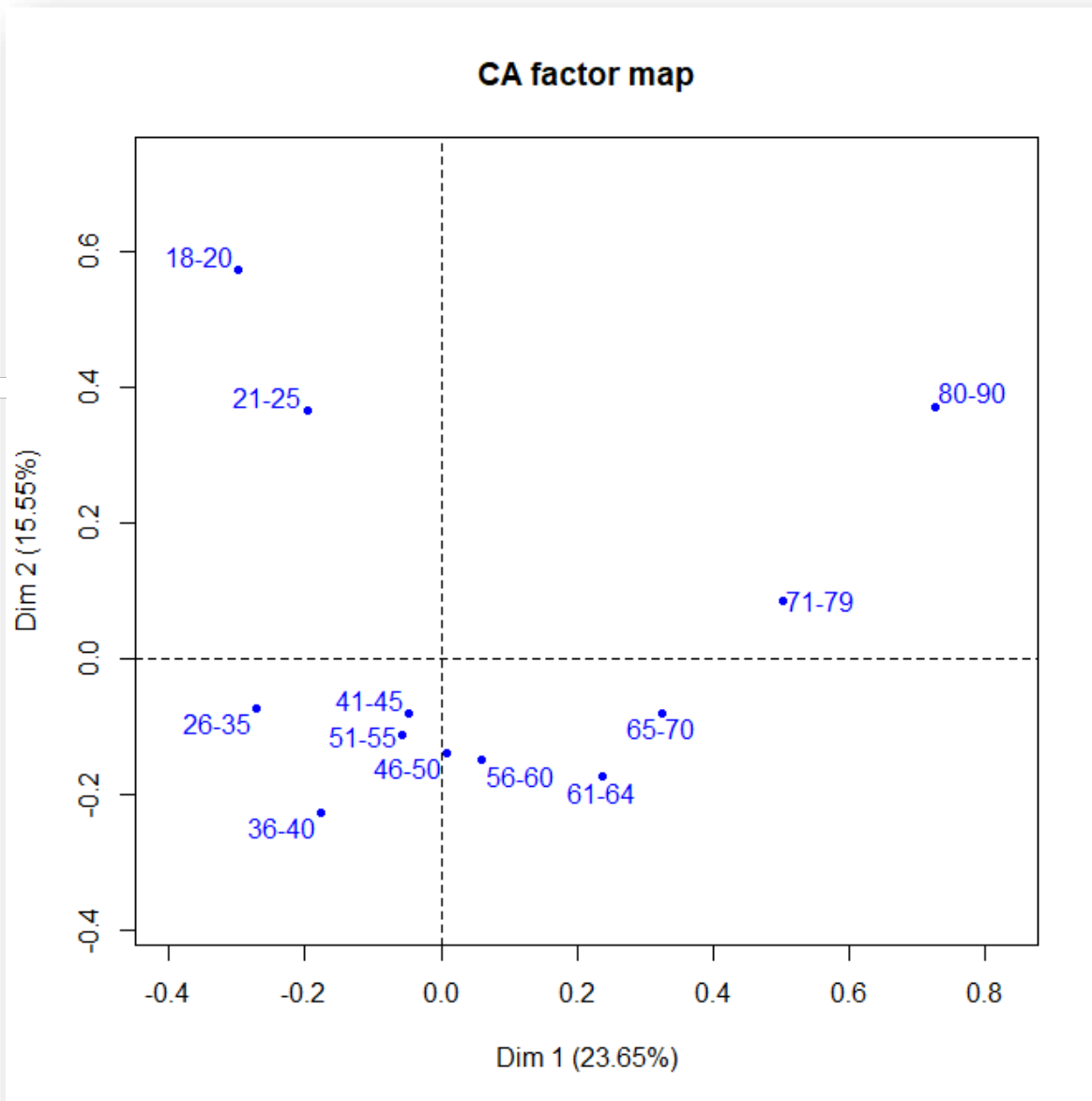
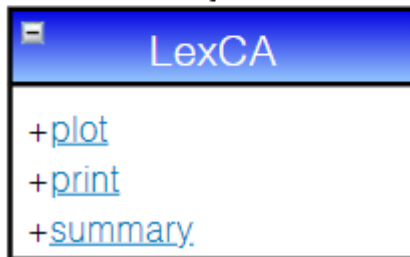
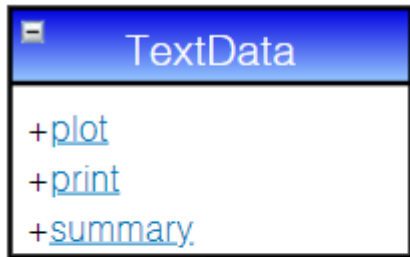


peace in the world: 6

LexCA function (Correspondence Analysis) (CA)



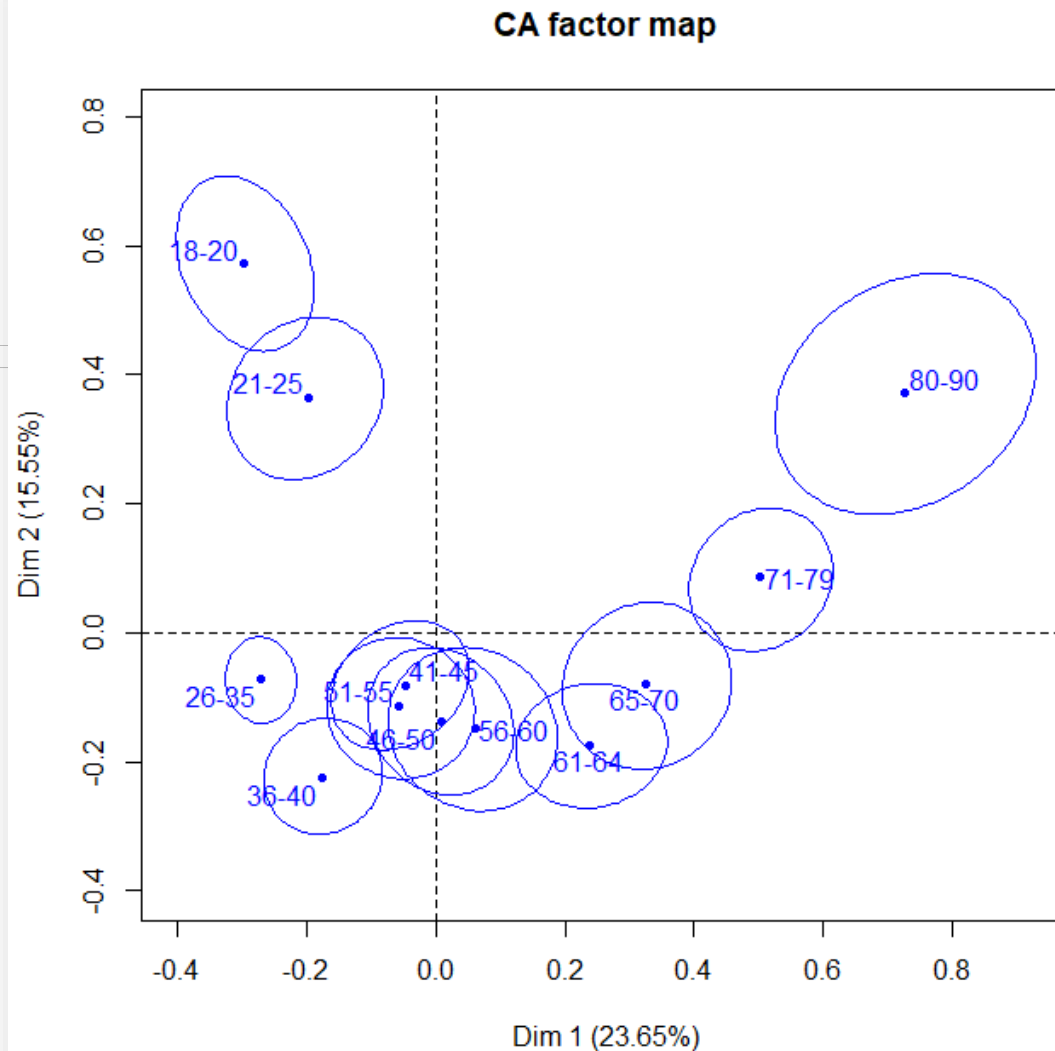
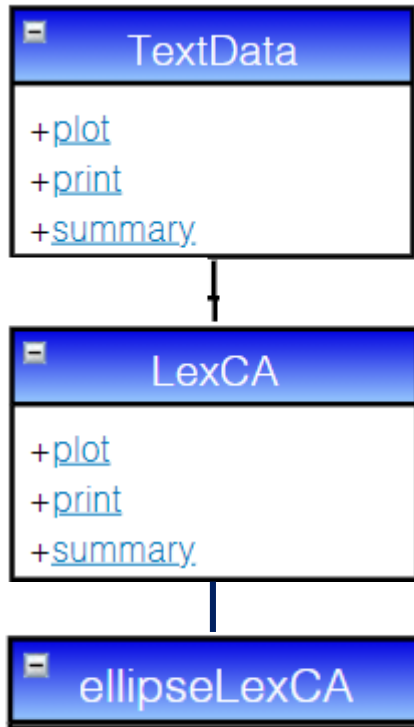
```
res.LexCA.2F <- LexCA(res.TD, ncp=2, graph=FALSE)  
plot(res.LexCA.2F, selWord=NULL)
```



LexCA function (Correspondence Analysis) (CA) / ellipseLexCA

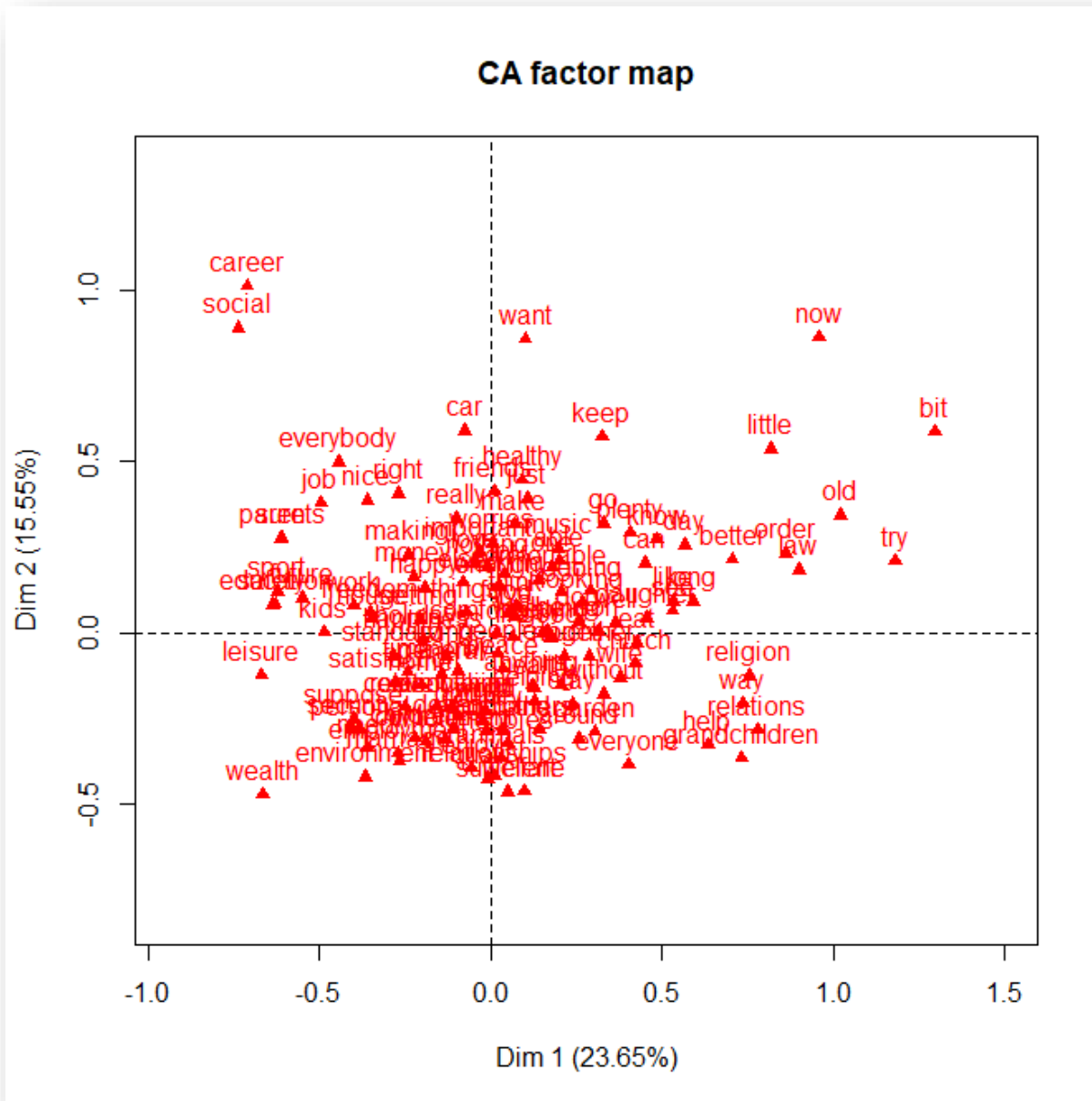


```
res.LexCA.2F <- LexCA(res.TD, ncp=2, graph=FALSE)  
ellipseLexCA(res.LexCA.2F, selWord=NULL, selDoc="ALL")
```



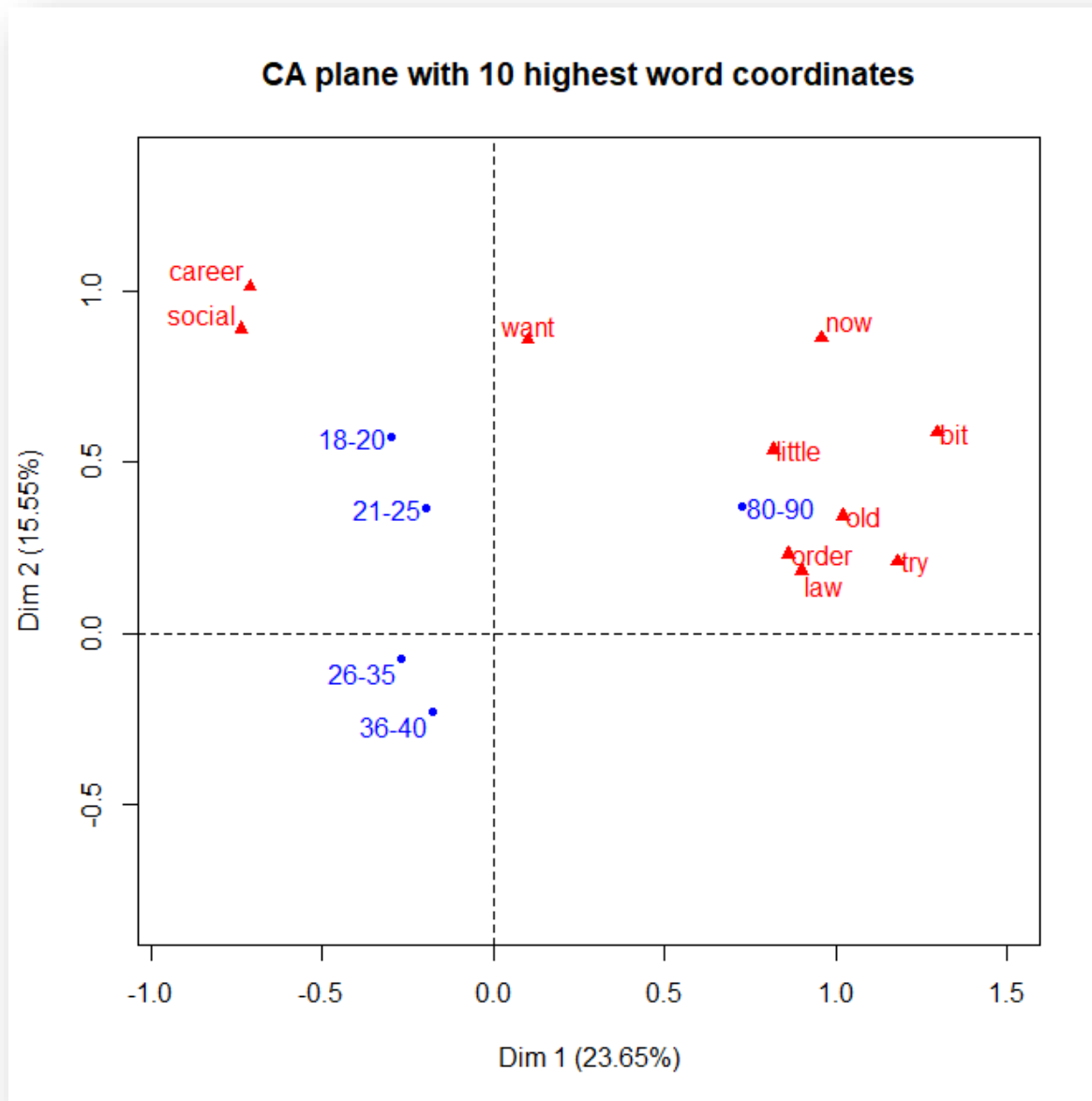


plot(res.LexCA.2F, selDoc=NULL)



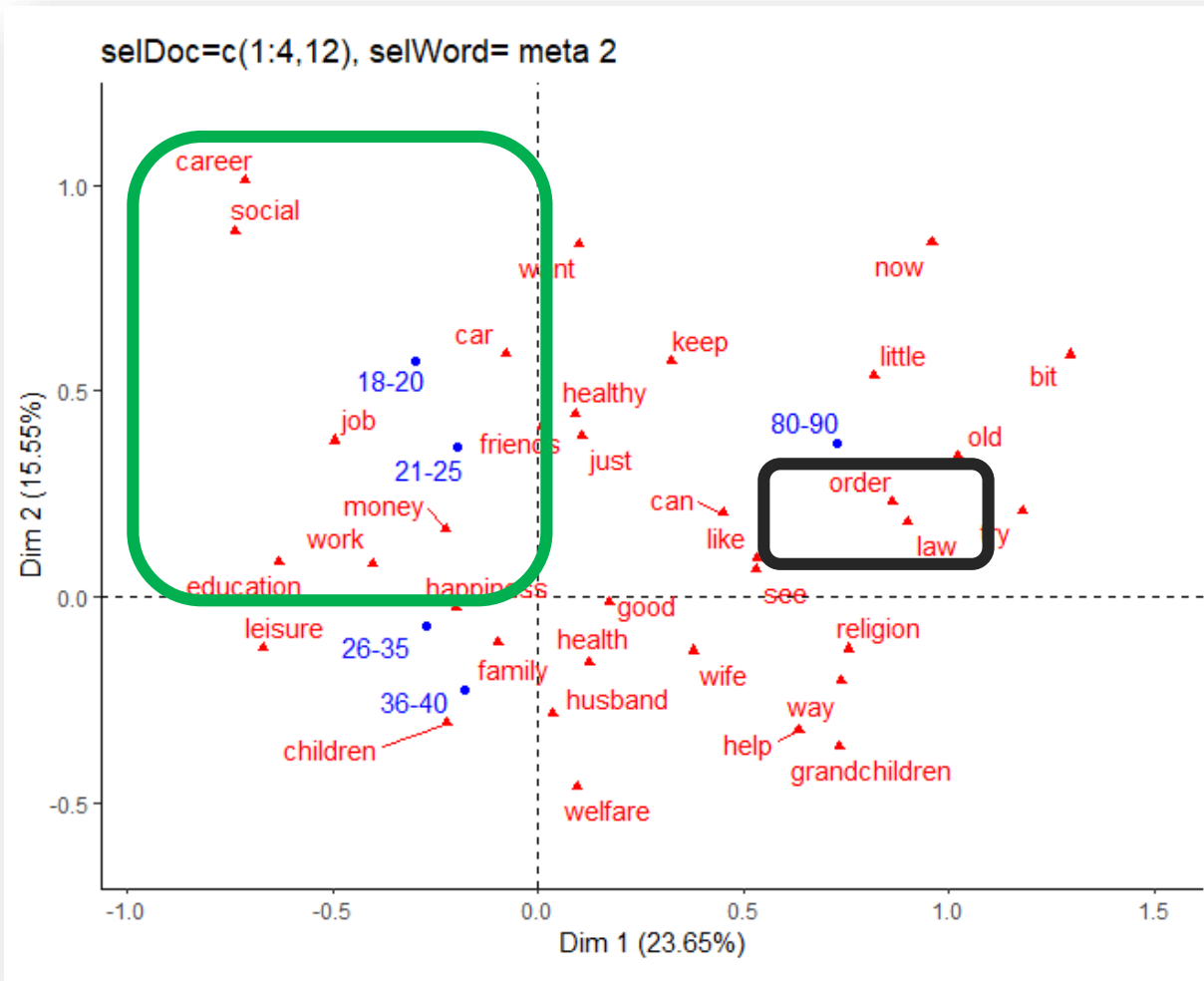


```
plot(res.LexCA.2F, selDoc=c(1,2,3,4, 12), selWord="coord 10")
```





```
plot(res.LexCA.2F, selDoc=c(1:4,12), selWord="meta 2")
```



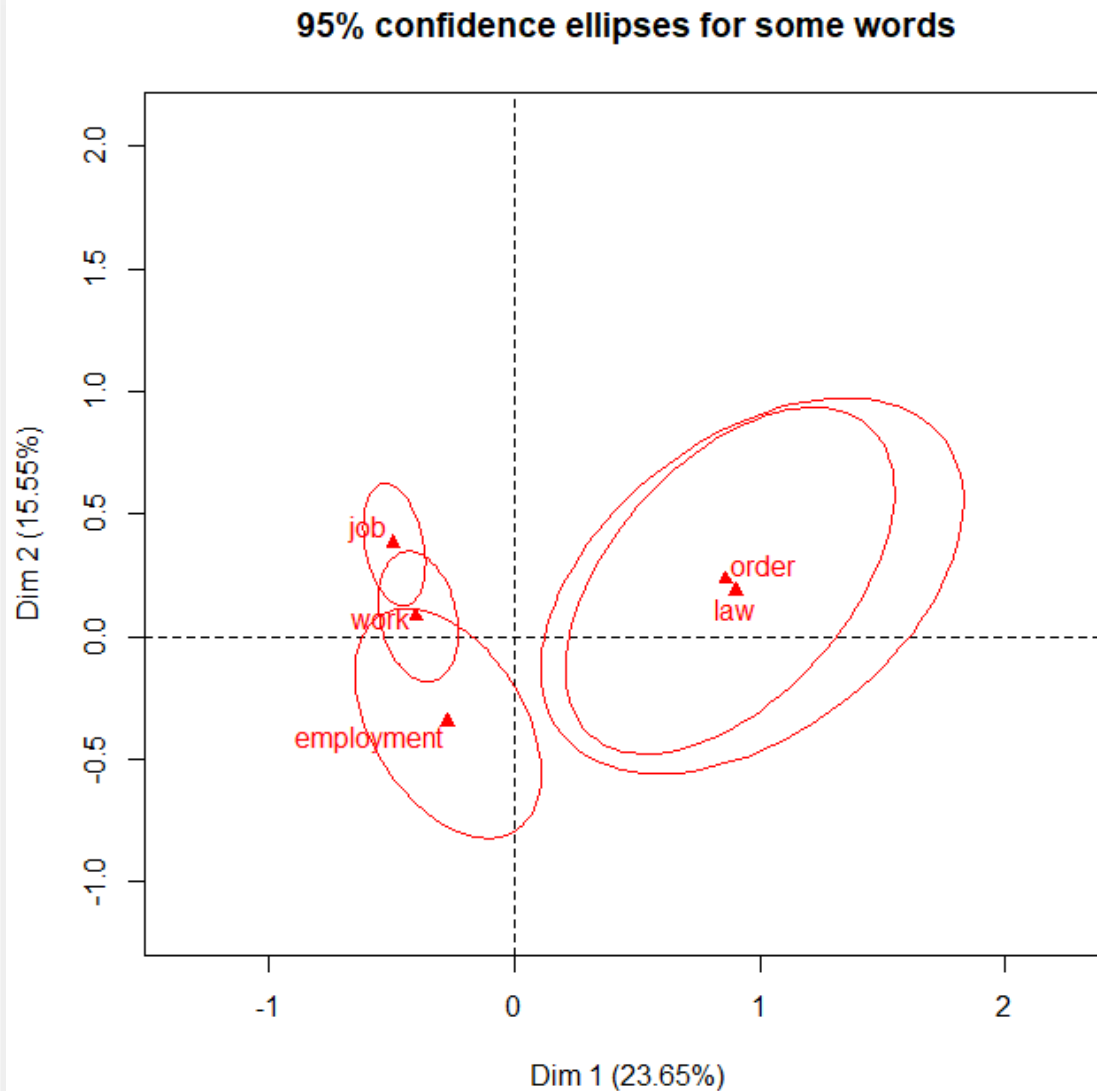
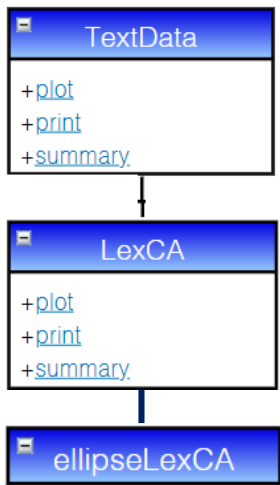
cos2

meta 2:

the words with a contribution over 2 times the average word contribution on any of the two axes are plotted.



```
ellipseLexCA(ellipseLexCA(res.LexCA.2F, selDoc=NULL,  
selWord=c("car","order","law", "job", "work", "employment", "money")))
```



LexHCca function

Hierarchical Clustering of Documents/Words on Textual Correspondence Analysis Coordinates)



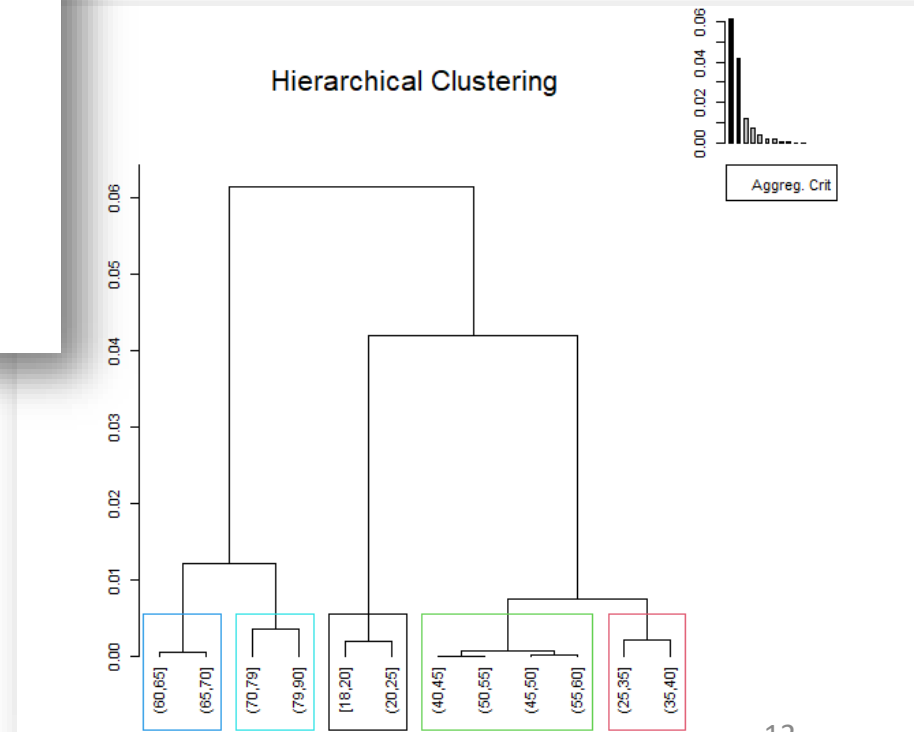
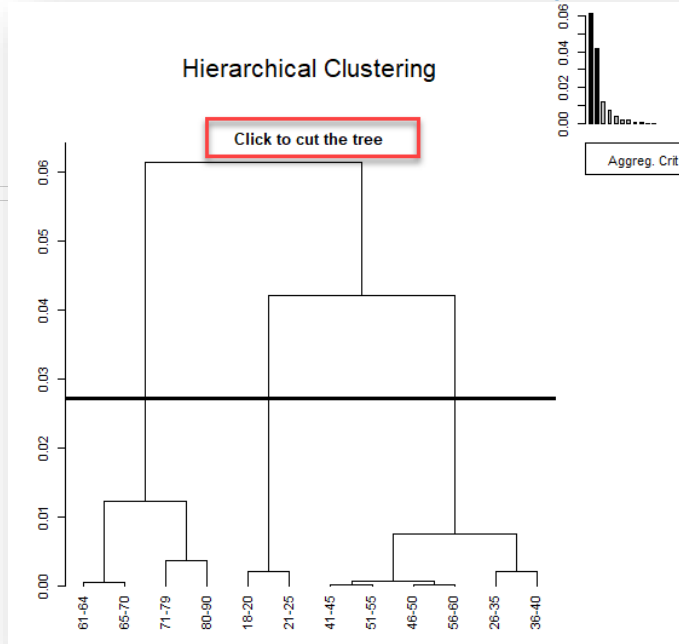
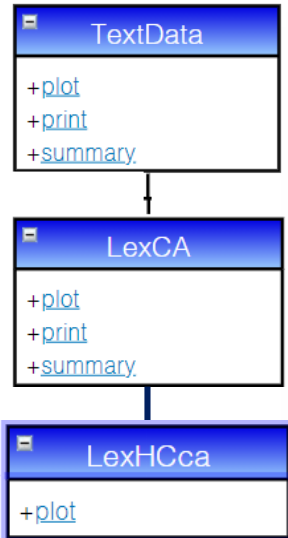
```
res.LexCA.2F <- LexCA(res.TD, ncp=2, graph=FALSE)
```

```
res.LexHCca.2F.5C <- LexHCca(res.LexCA.2F, nb.clust="click")
```

```
res.LexHCca.2F.5C <- LexHCca(res.LexCA.2F, nb.clust=5)
```

```
cluster.CA="docs"
```

```
cluster.CA="words"
```

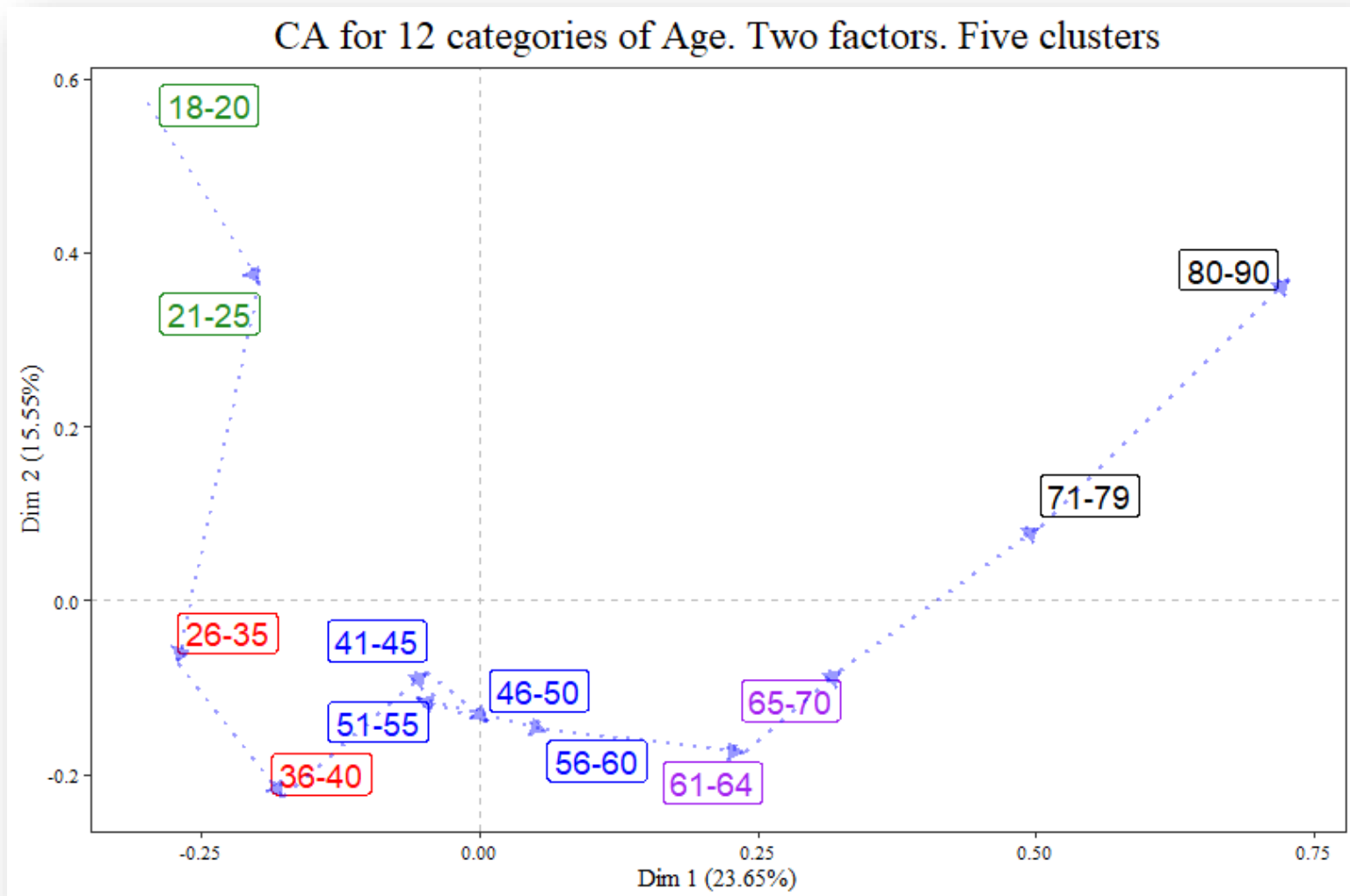


```
class(res.LexHCca.2F.5C$call$t$tree)
```

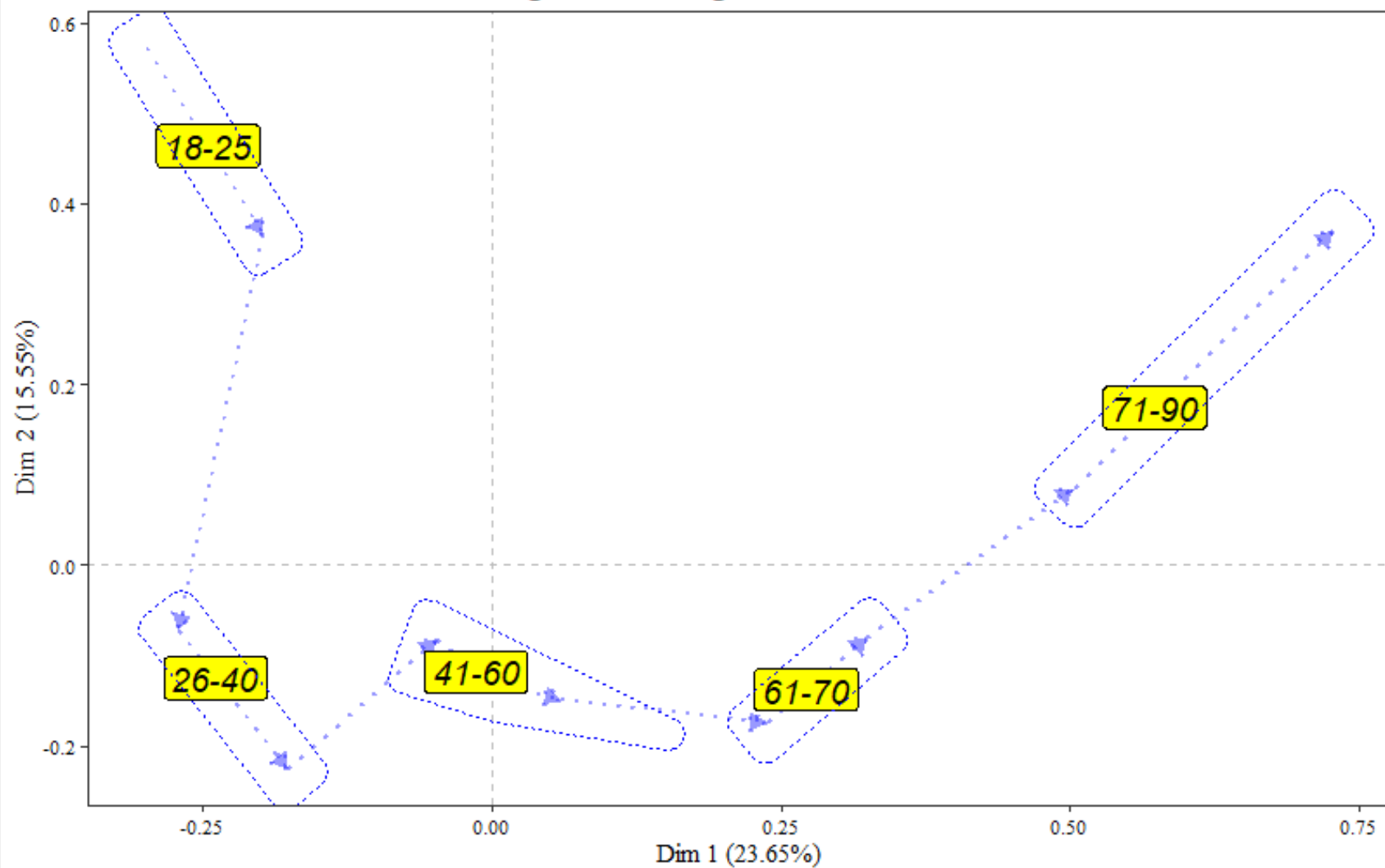
```
hclust
```



```
c.palette <- c("forestgreen", "red", "blue", "purple", "black")  
plot(res.LexHCca.2F.5C, plot=c("labels", "traject"),  
     labels=c(size=5, force=0.2, set.seed=998, rect=TRUE), palette=c.palette,  
     traject=c(color="blue", linetype="dotted", alpha.t=.4),  
     title=c(text="CA for 12 categories of Age. Two factors. Five clusters"))
```

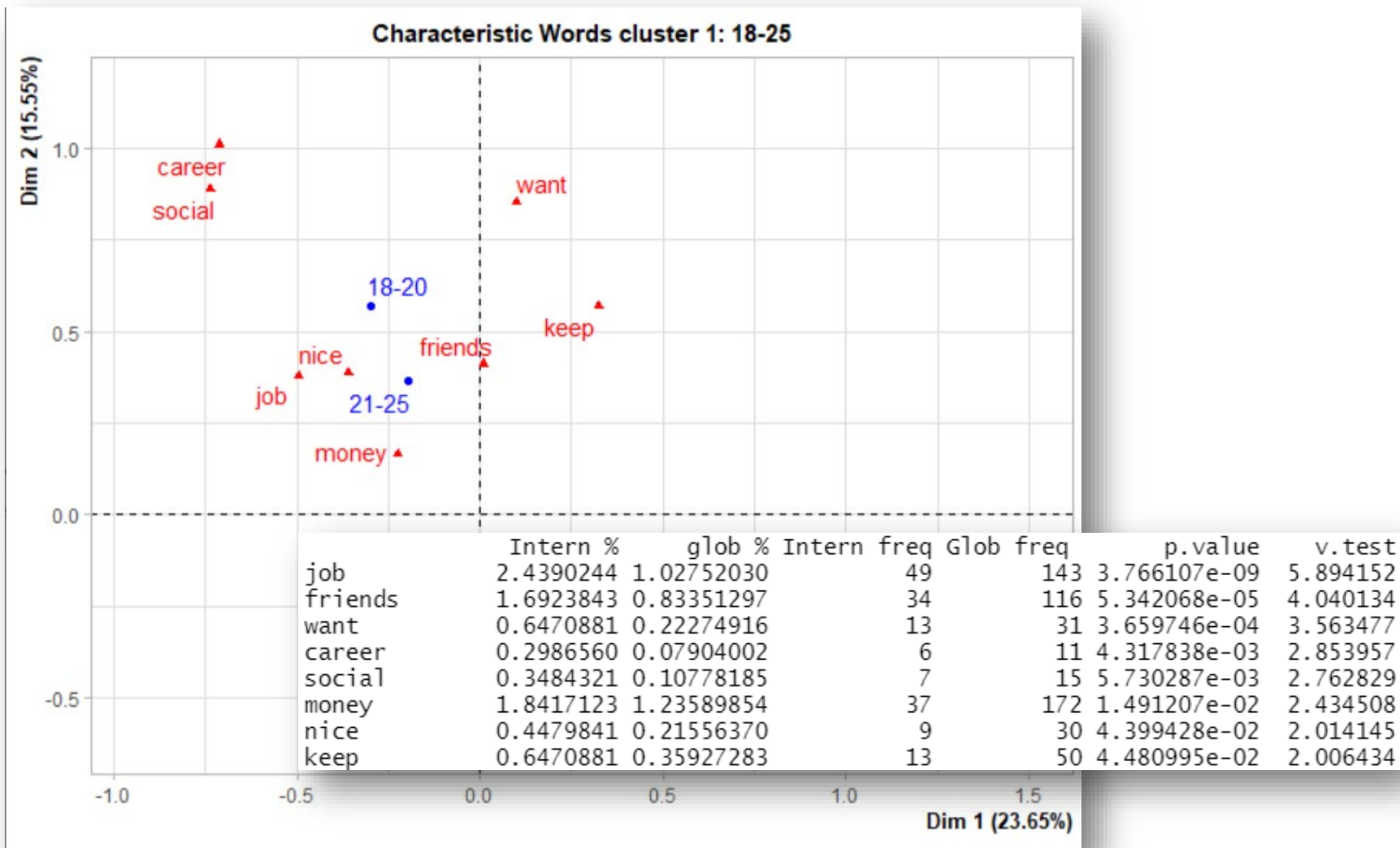


CA for 12 categories of Age. Two factors. Five clusters





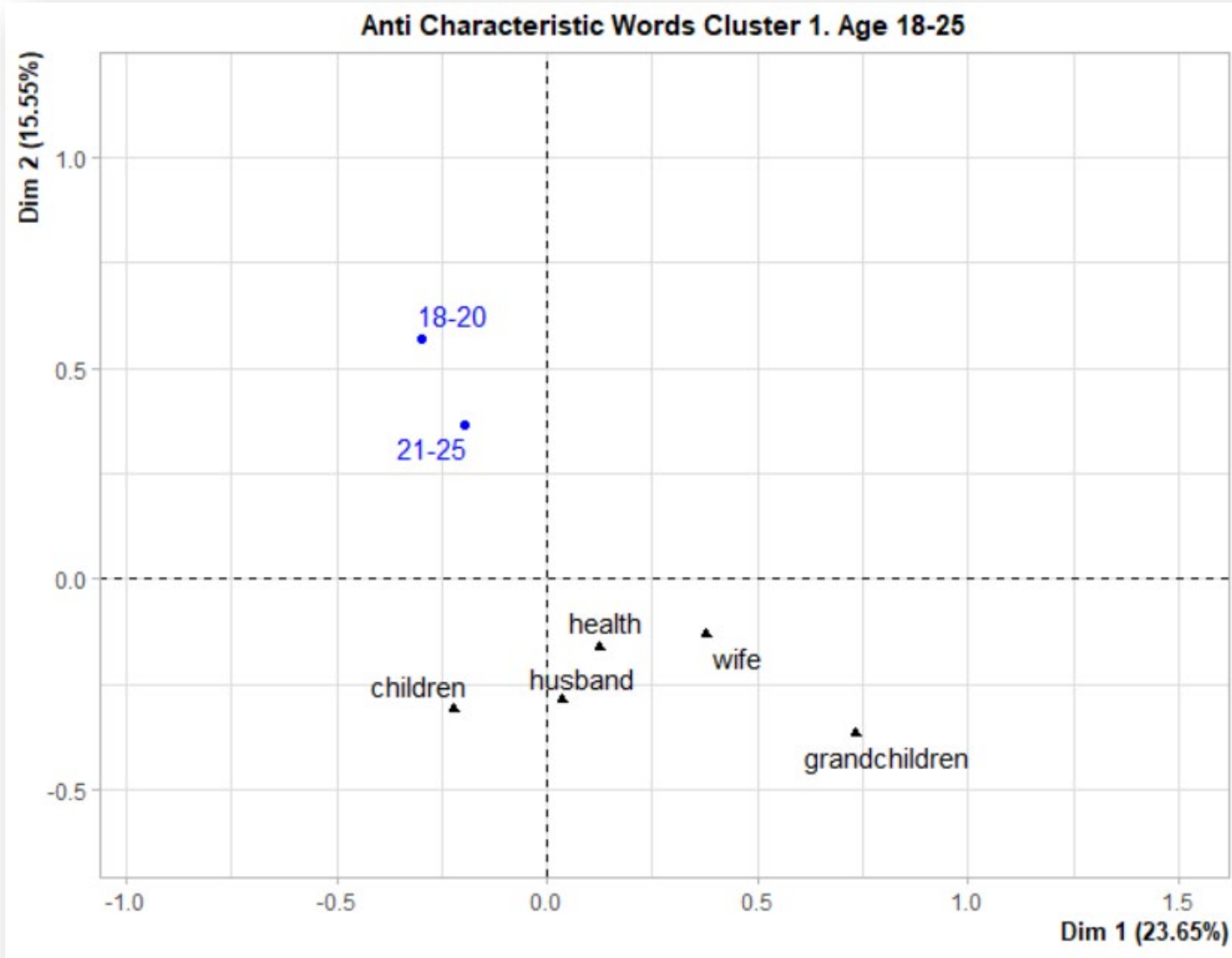
```
df <- as.data.frame(res.LexHCca.2F.5C$description$desc.cluster.doc$words$cluster_1)
words_1 <- rownames(df[df$v.test>0,])
plot(res.LexCA.2F, selDoc=c(1,2), selWord=words_1, graph.type = "ggplot",
     title="Characteristic Words 18-25")
```





```
anti_words_1 <- rownames(df[df$v.test<0,])
```

```
plot(res.LexCA.2F, selDoc=c(1,2), selWord=anti_words_1, col.word = "black", graph.type = "ggplot",  
title="Anti Characteristic Words 18-25")
```

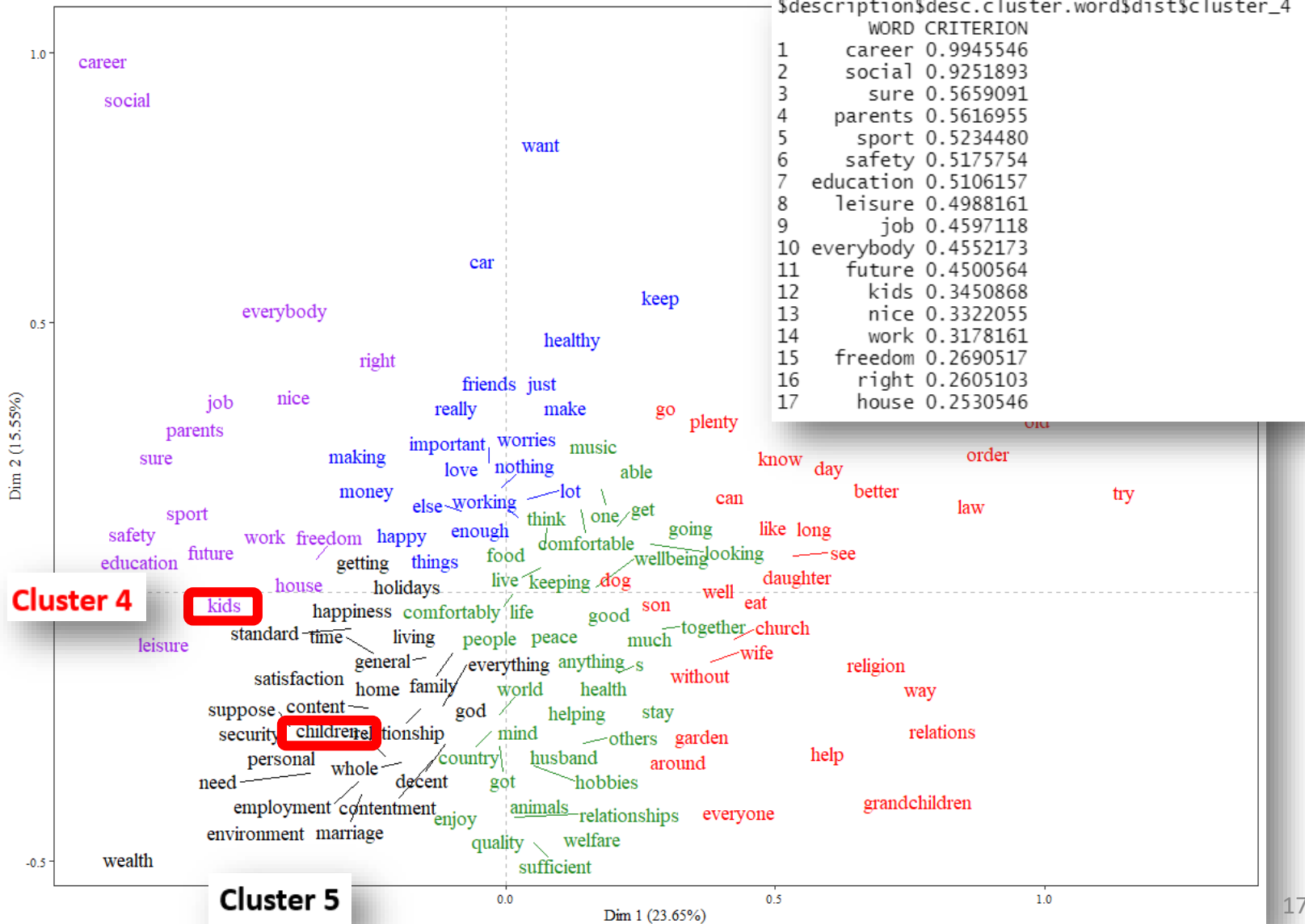


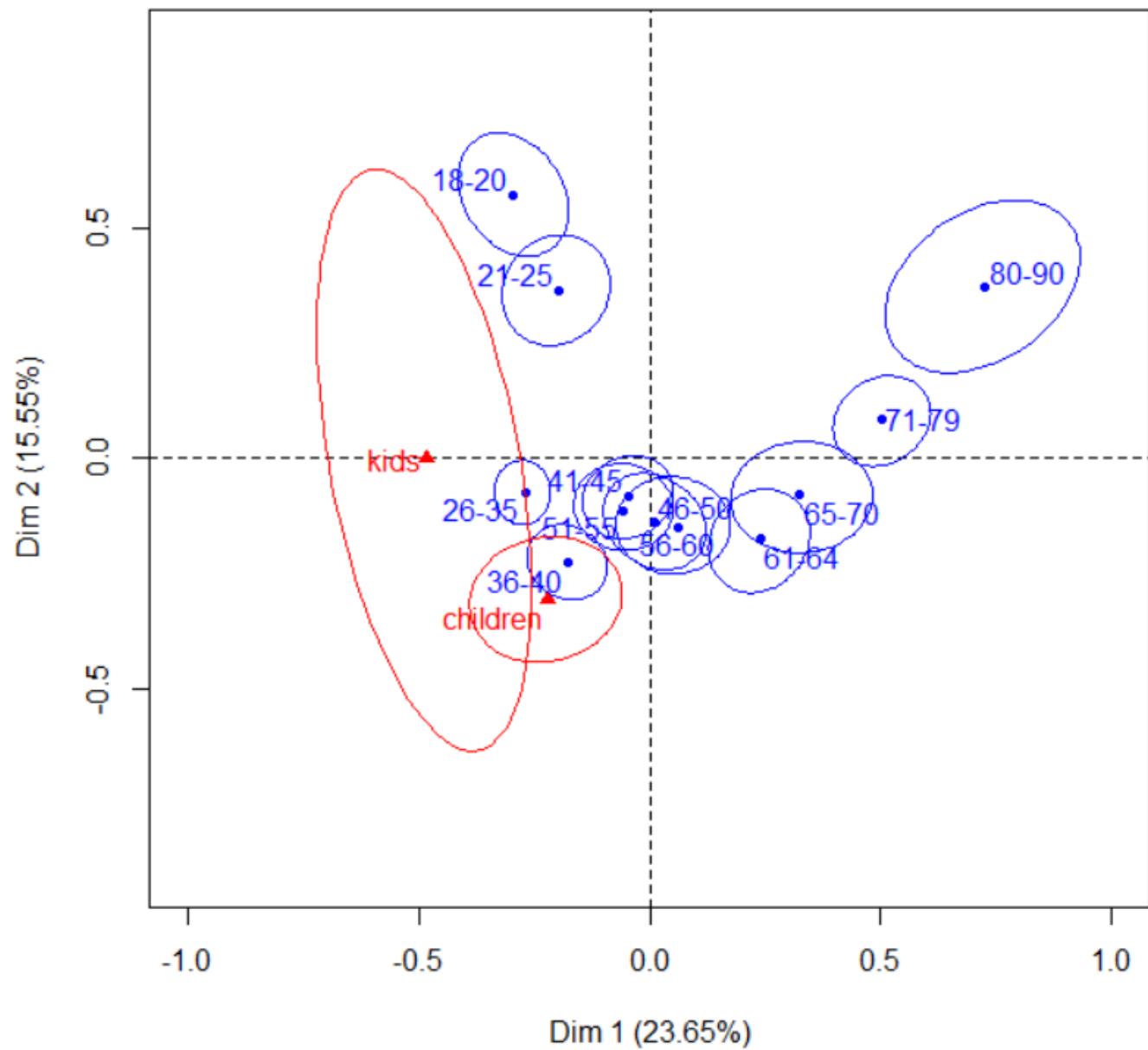
LexHCca function Hierarchical Clustering of Words

Words.C <- LexHCca(res.LexCA.2F, cluster.CA="words", ...

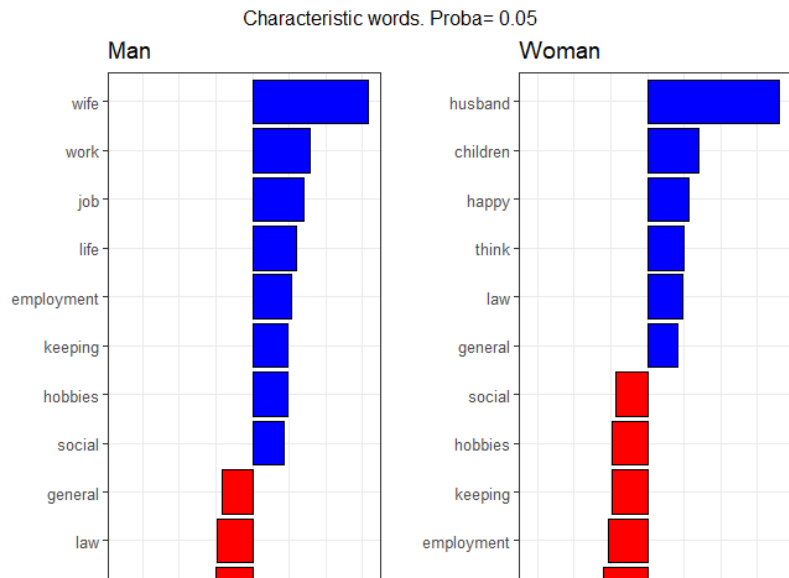


CA for 12 categories of Age. Two factors. Five clusters





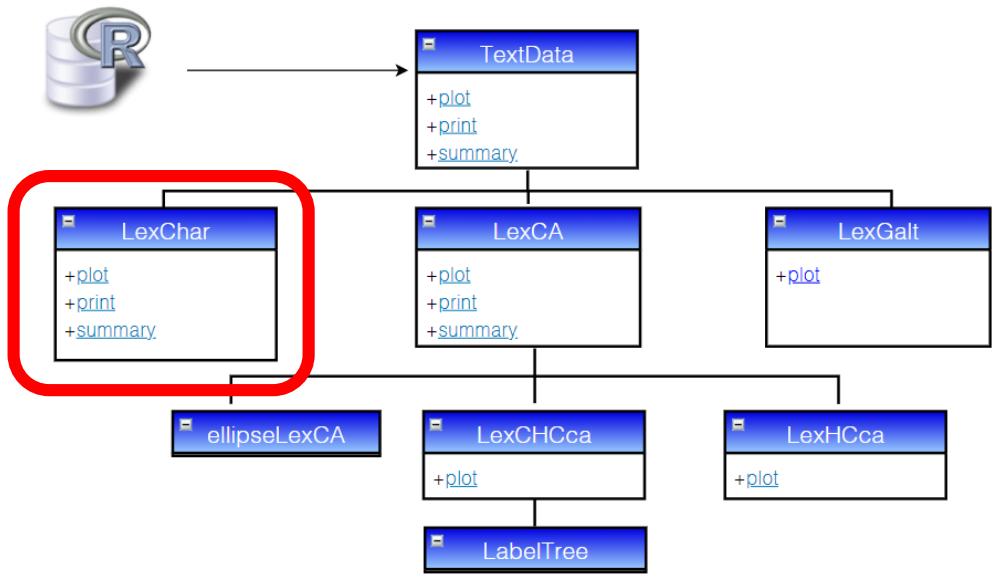
Another analysis not presented here



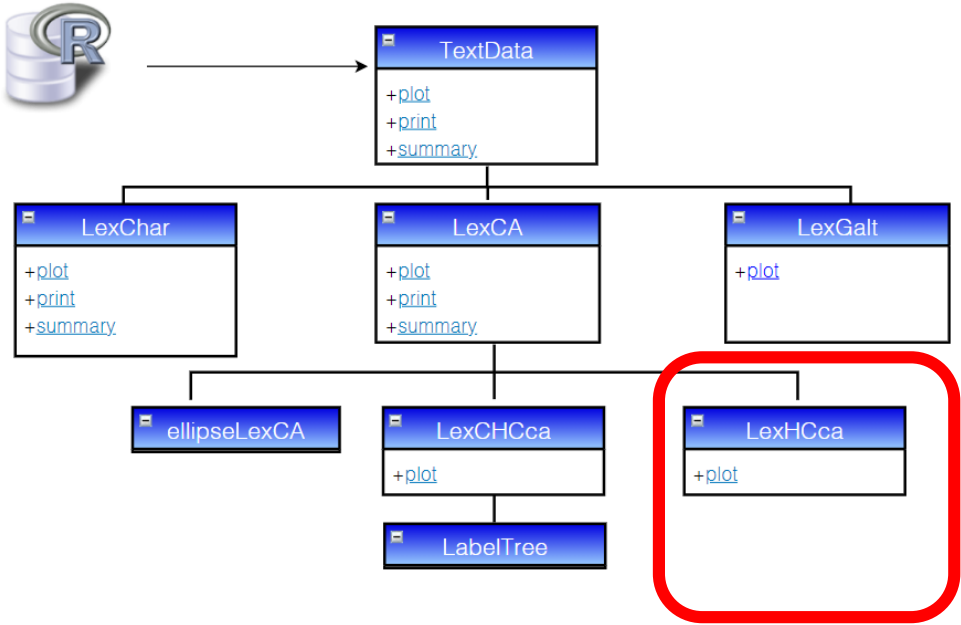
\$woman	DOCUMENT	CRITERION
1	4	0
2	5	0
3	7	0
4	8	0
5	9	0
6	10	0
7	12	0
8	13	0
9	15	0
10	16	0

to be healthy.just to live long enough to see the children grow up, I do not think there is a lot t

-----TEXT-----
 health. happiness, money, family
 to be happy. healthy, have enough to eat, enough money to live on
 health. happiness
 health,. keeping going, family, going out, shopping, visiting
 husband. new baby grand daughter, life in general,
 good health. happiness, togetherness,
 family. friends, pets,
 my family really. health, walking
 my children.my husband, my family and relations, health and wellbeing of family

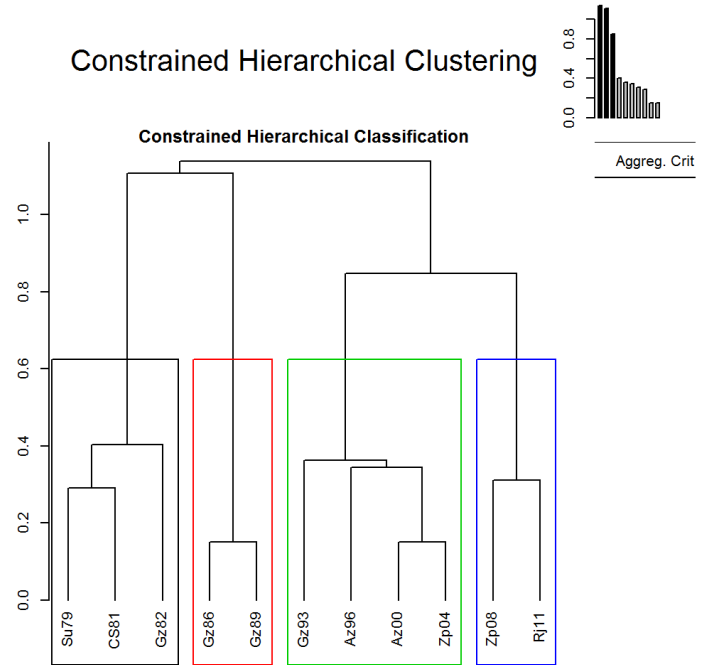


Another analysis not presented here

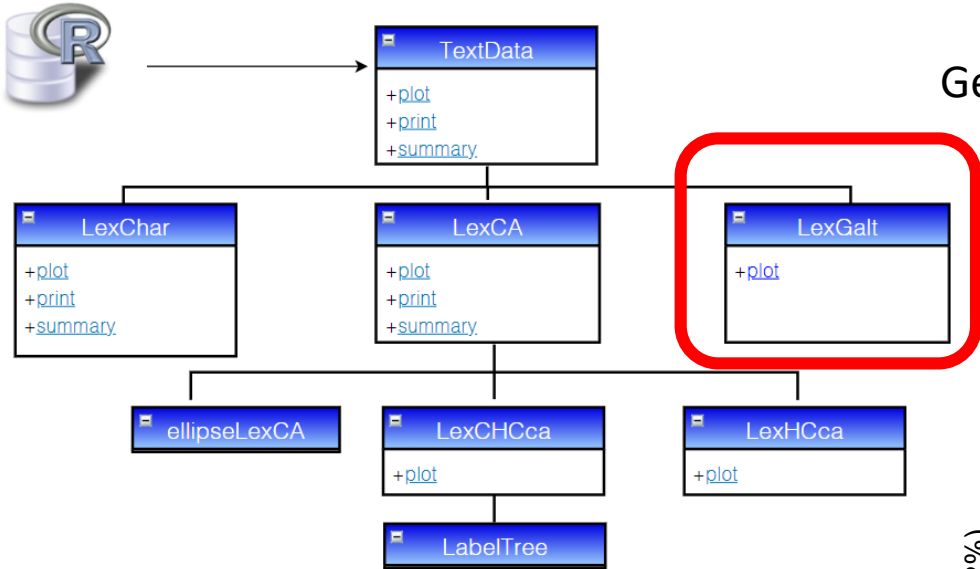


Chronologically Constrained Cluster

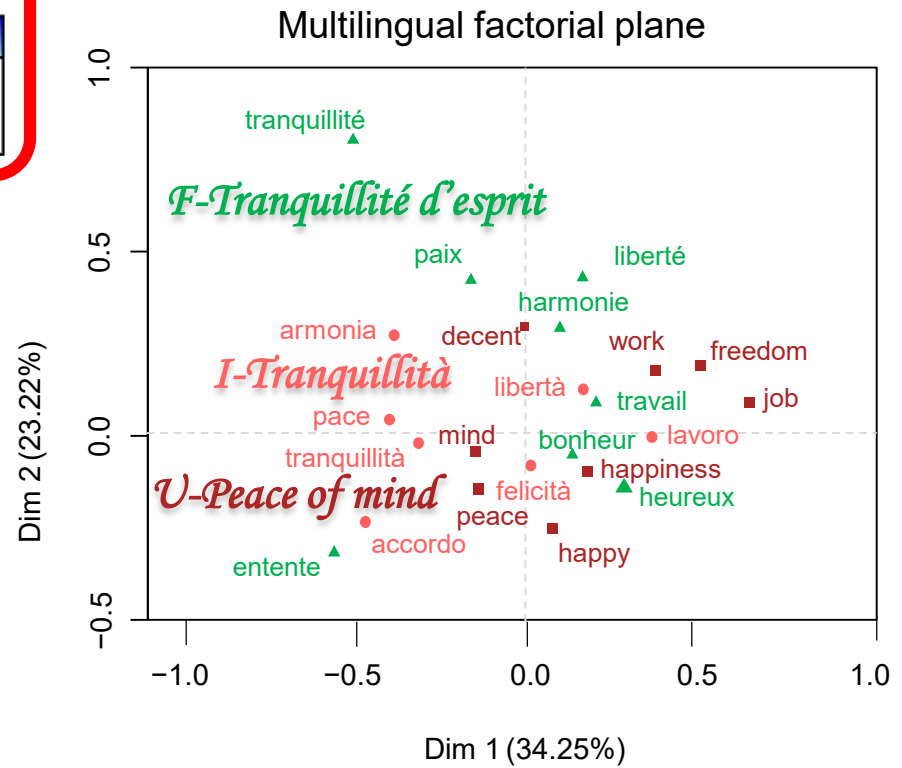
Constrained Hierarchical Clustering



Another analysis not presented here



Generalised Aggregate Lexical Table (LexGalt)



5 May 2021

Las Médulas. Spanish Roman over ground gold mine

Thank you very much

Authors:

Ramón Álvarez-Esteban (2) (Maintainer)

Mónica Bécue-Bertaut (1)

Josep-Anton Sánchez-Espigares (1)

Belchin Kostov (1)

(1) UPC Universitat Politècnica de Catalunya / Spain

(2) University of Leon / Spain. ramon.alvarez@unileon.es

ramon.alvarez@unileon.es