

PCA on Finnish Municipalities

Dimensionality reduction with PCA on Finnish municipality demographics

Niko Miller

10.04.2022

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Objective and research question	1
1.3	Scope	1
2	Data and Exploratory Analysis	1
2.1	Variable Descriptions	1
2.2	Univariate Analysis	2
2.2.1	Descriptive Statistics	2
2.2.2	Distribution Plots	3
2.3	Bivariate Analysis	5
2.3.1	Pearson's Correlation	5
3	Principal Component Analysis (PCA)	6
3.1	Scree Plot	6
3.2	Principal Components 1 and 2	7
3.3	Principal Components 3 and 4	9
4	Discussion and Conclusions	10

1 Introduction

1.1 Motivation

Tilastokeskus provides data on 32 demographic variables of Finnish municipalities [Statistics Finland, 2022]. Those variables are in many cases correlated and some are perhaps better in explaining differences between municipalities. Hence, an interesting question is whether we could reduce dimensionality of the data and describe municipalities using only a few variables.

1.2 Objective and research question

1.3 Scope

2 Data and Exploratory Analysis

I used PXWEB API [Magnusson et al., 2019] to retrieve a data set on key ratios for all of Finland's municipalities during 1987-2019 according to the 2020 municipality classification. The data set is provided by Tilastokeskus [Statistics Finland, 2022].

2.1 Variable Descriptions

Table 1. Variable Descriptions

Original	Modified	Explanation
Taajama-aste, %	Prop.Urban.Areas	NA
Väkiluku	Pop	NA
Väkiluvun muutos edellisestä vuodesta, %	Pop.Change	NA
Alle 15-vuotiaiden osuus väestöstä, %	Prop.Below15	NA
15-64 -vuotiaiden osuus väestöstä, %	Prop.15to64	NA
Yli 64-vuotiaiden osuus väestöstä, %	Prop.Over64	NA
Ruotsinkielisten osuus väestöstä, %	Prop.Swedish	NA
Ulkomaan kansalaisten osuus väestöstä, %	Prop.Foreign	NA
Syntyneiden enemmyys, henkilöä	Pop.Growth	NA
Kuntien välinen muuttovoitto/-tappio, henkilöä	Migr.Gain	NA
Perheiden lukumäärä	Families	NA
Asuntokuntien lukumäärä	Households	NA
Rivi- ja pientaloissa asuvien asuntokuntien osuus, %	Prop.Households.rowSmall	NA
Vuokra-asunnoissa asuvien asuntokuntien osuus, %	Prop.Households.Rent	NA
Vähintään toisen asteen tutkinnon suorittaneiden osuus 15 vuotta täyttäneistä, %	Prop.Educ.Degree2	NA
Korkea-asteen tutkinnon suorittaneiden osuus 15 vuotta täyttäneistä, %	Prop.Educ.High	NA
Alueella asuvan työllisen työvoiman määrä	Employed	NA
Työllisyysaste, %	Empl.Rate	NA
Asuinkunnassaan työssäkäyvien osuus, %	Prop.Households.Empl	NA
Työttömien osuus työvoimasta, %	Prop.Unempl	NA
Eläkeläisten osuus väestöstä, %	Prop.Pension	NA
Taloudellinen huoltosuhde	Support.Ratio	NA
Alueella olevien työpaikkojen lukumäärä	Jobs	NA
Alkutuotannon työpaikkojen osuus, %	Prop.Primary.Sector	NA
Jalostuksen työpaikkojen osuus, %	Prop.Secondary.Sector	NA
Palvelujen työpaikkojen osuus, %	Prop.Services.Sector	NA
Työpaikkaomavaraisuus	Jobs.Self.Suff	NA
Vuosikate, euroa/asukas	Margin.Citizen	NA
Lainakanta, euroa/asukas	Loan.Citizen	NA
Konsernin lainakanta, euroa/asukas	Concern.Loan.Citizen	NA
Opetus- ja kulttuuritoiminta yhteensä, nettokäyttökustannukset, euroa/asukas	Educ.Cult.Citizen	NA
Sosiaali- ja terveystoiminta yhteensä, nettokäyttökustannukset, euroa/asukas	Soc.Health.Citizen	NA

2.2 Univariate Analysis

2.2.1 Descriptive Statistics

Table 2. Descriptive Statistics

Var	Min	1st Qu.	Median	Mean	3rd Qu.	Max	Std.
Prop.Urban.Areas	0.0	46.9	60.2	61.4	77.1	100.0	22.2
Pop	92.0	2848.2	6298.5	17845.9	15259.5	643272.0	47989.0
Pop.Change	-5.2	-1.7	-1.0	-0.8	0.0	3.5	1.4
Prop.Below15	1.1	13.1	15.1	15.6	17.9	33.6	4.0
Prop.15to64	48.9	55.0	57.3	57.6	59.9	68.9	3.5
Prop.Over64	9.6	22.2	26.6	26.8	31.6	42.4	6.3
Prop.Swedish	0.0	0.1	0.3	10.9	1.0	93.5	26.4
Prop.Foreign	0.3	1.2	1.9	2.8	3.1	17.3	2.7
Pop.Growth	-470.0	-51.5	-23.0	-10.4	0.0	1649.0	170.4
Migr.Gain	-763.0	-73.2	-33.0	-1.5	-7.8	5027.0	376.4
Families	22.0	756.0	1745.5	4768.0	4223.5	158063.0	12139.4
Households	54.0	1335.8	2861.5	8659.5	6917.8	330933.0	24255.1
Prop.Households.rowSmall	13.3	74.9	88.3	81.9	93.8	98.8	16.2
Prop.Households.Rent	6.8	16.4	19.6	20.8	23.3	49.2	7.1
Prop.Educ.Degree2	55.7	64.1	67.1	67.3	70.4	82.2	4.9
Prop.Educ.High	11.1	18.7	22.0	23.1	26.6	57.6	6.4
Employed	44.0	1018.8	2412.0	7550.8	6150.5	309685.0	22305.5
Empl.Rate	56.1	66.5	70.6	70.8	75.1	87.2	6.3
Prop.Households.Empl	19.6	43.6	60.6	58.7	74.1	92.1	18.3
Prop.Unempl	1.6	8.1	10.8	10.9	13.3	20.7	4.0
Prop.Pension	12.4	26.4	32.0	32.0	37.7	48.6	7.5
Support.Ratio	87.5	134.2	156.9	160.0	182.5	239.6	32.9
Jobs	24.0	866.5	2057.0	7559.6	5073.0	397346.0	26845.9
Prop.Primary.Sector	0.1	4.1	10.1	11.0	16.2	41.3	8.1
Prop.Secondary.Sector	2.5	18.0	23.5	24.3	30.1	63.0	10.0
Prop.Services.Sector	28.0	56.1	62.6	62.6	68.8	94.0	10.6
Jobs.Self.Suff	34.4	72.3	87.8	86.2	99.6	171.1	20.5
Margin.Citizen	-195.7	355.9	496.3	516.3	646.8	1590.2	251.0
Loan.Citizen	0.0	1871.6	2752.3	2792.9	3563.9	10508.4	1488.3
Concern.Loan.Citizen	0.0	3186.8	4332.8	4612.2	5869.9	12290.8	2374.6
Educ.Cult.Citizen	1043.5	1691.5	1847.6	1902.0	2063.8	3081.2	307.6
Soc.Health.Citizen	1094.7	3163.6	3610.4	3618.3	4157.6	6009.0	800.6

2.2.2 Distribution Plots

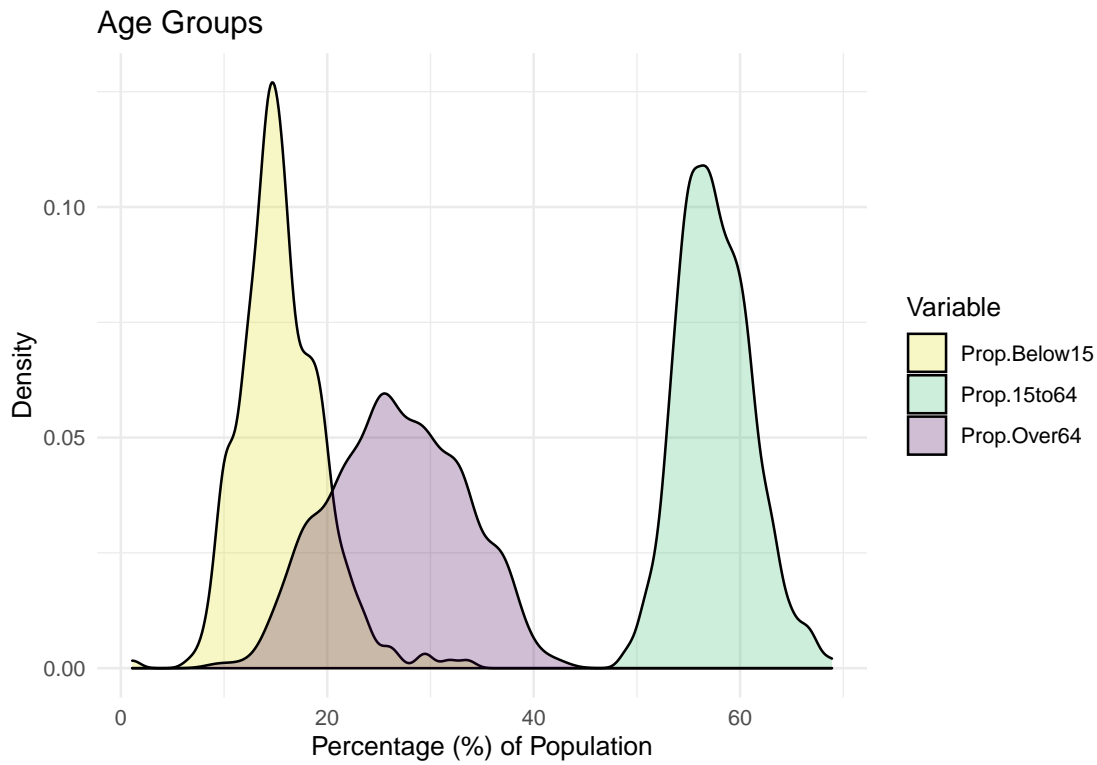


Figure 1. Probability density estimates for age groups

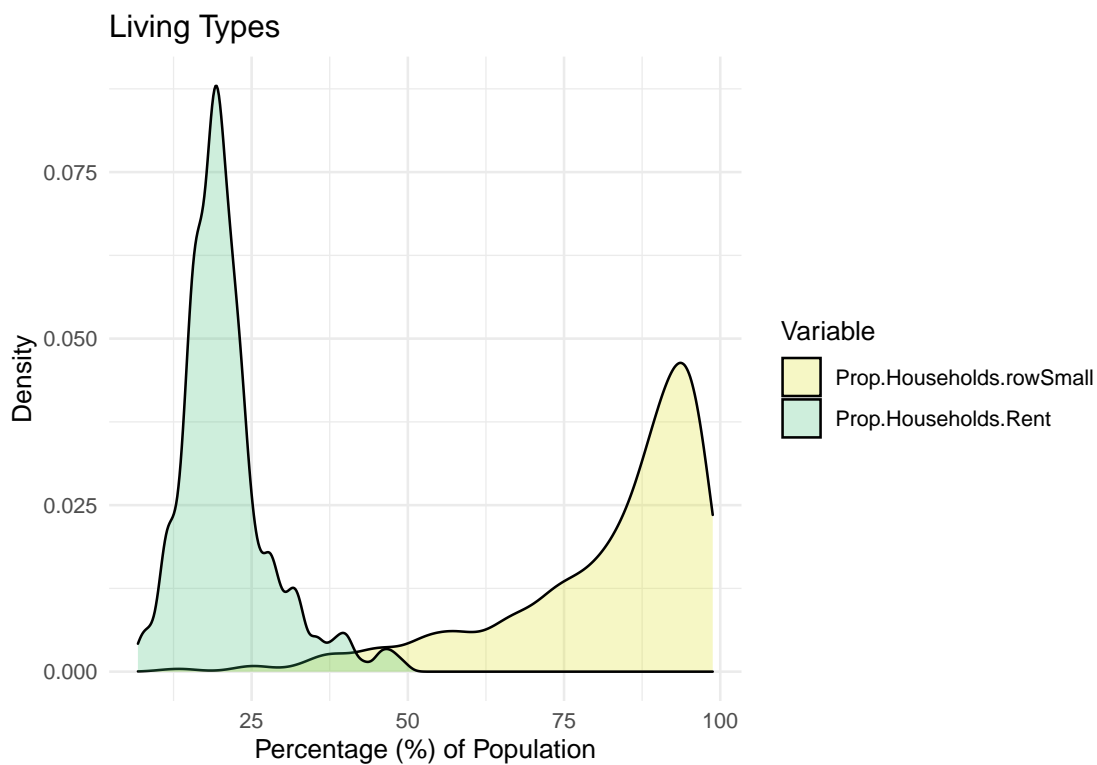


Figure 2. Probability density estimates for living types

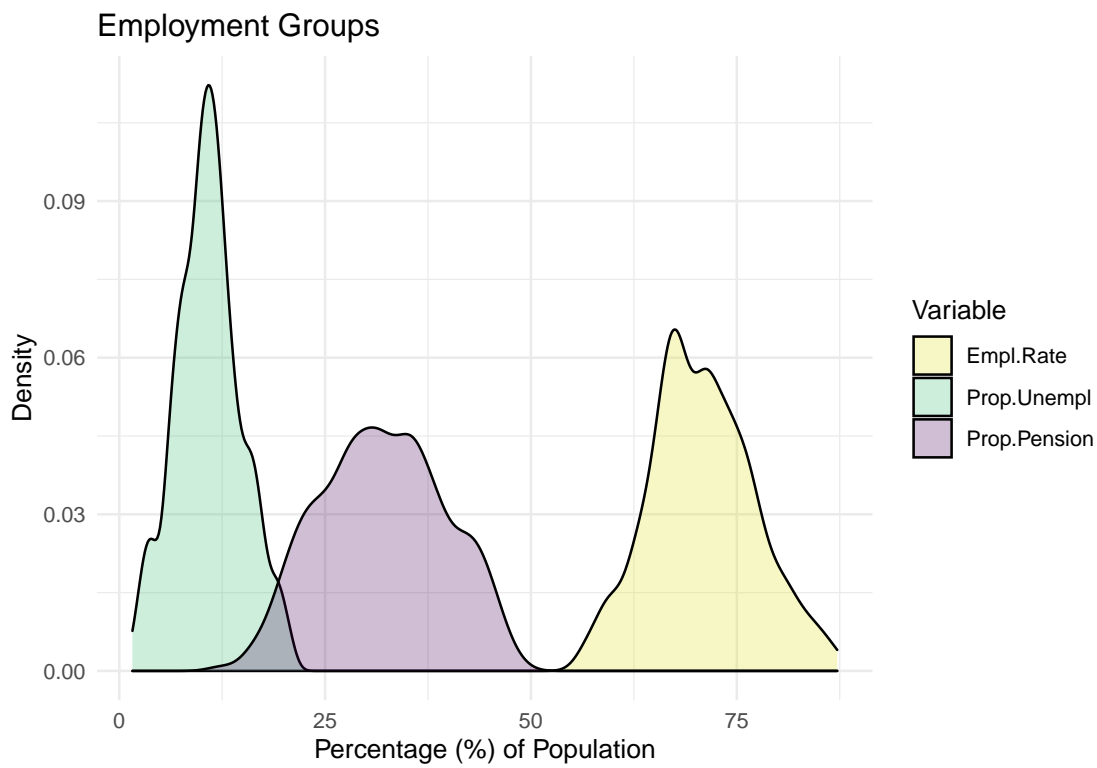


Figure 3. Probability density estimates for employment groups

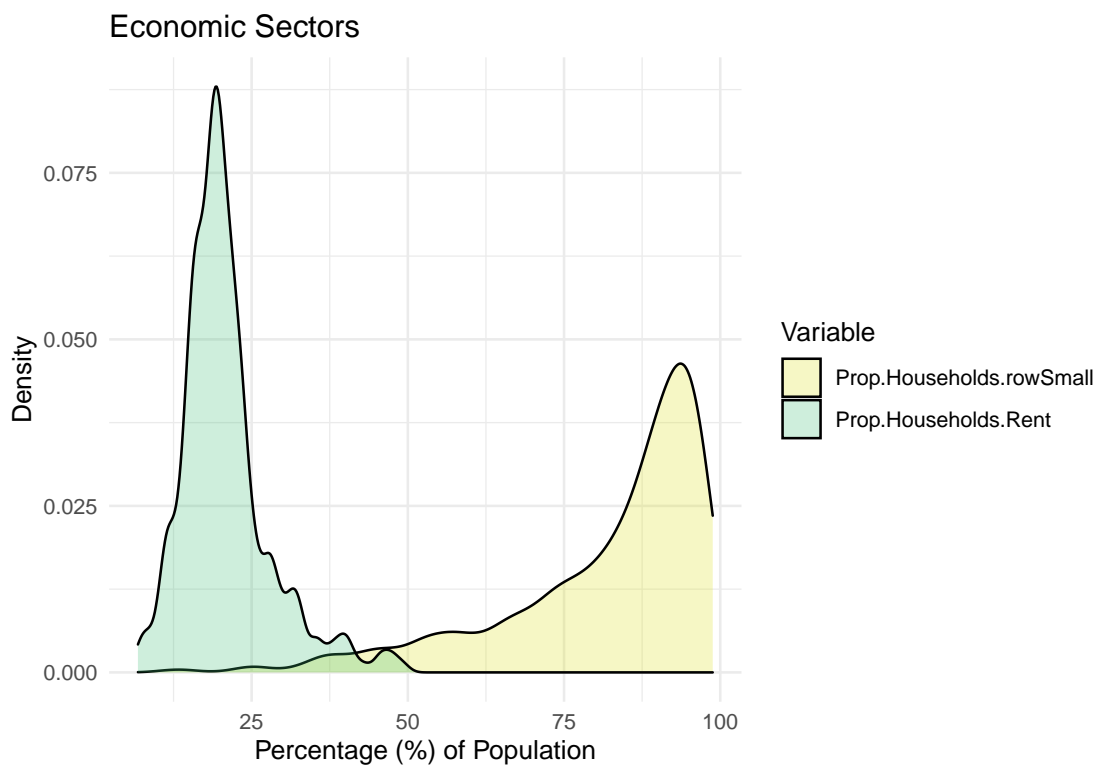


Figure 4. Probability density estimates for economic sectors

2.3 Bivariate Analysis

2.3.1 Pearson's Correlation

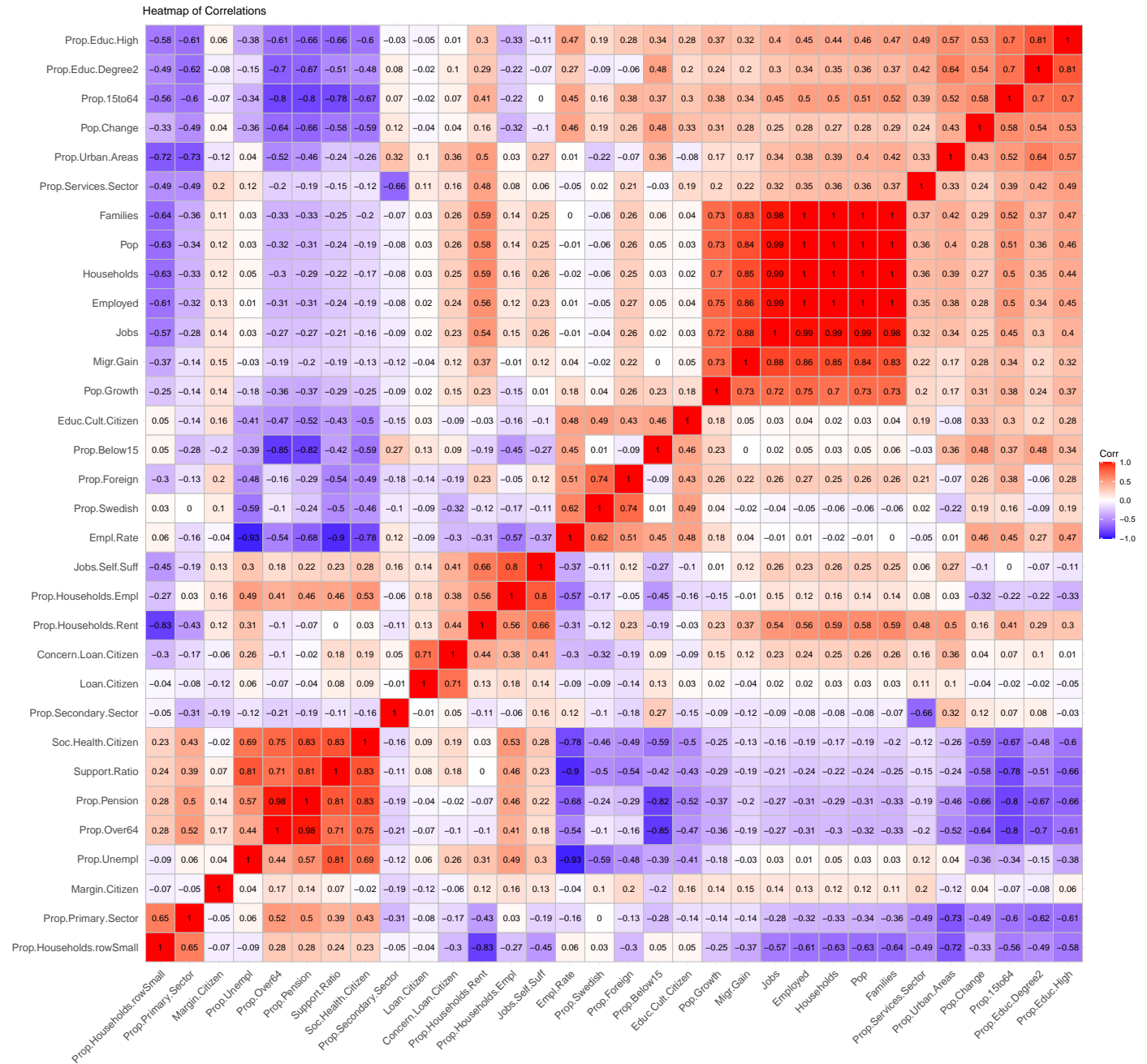


Figure 5. Heatmap of Pearson's correlation between all variables

3 Principal Component Analysis (PCA)

3.1 Scree Plot

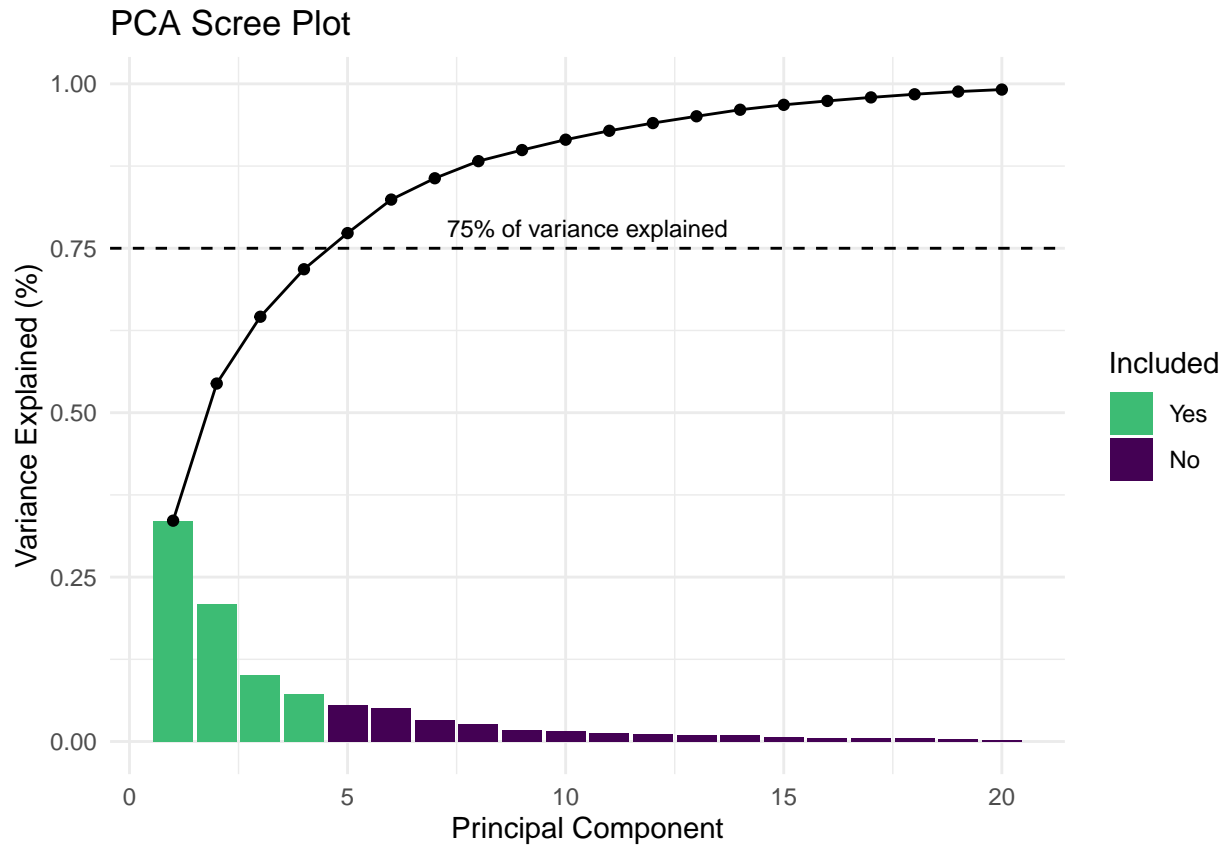


Figure 6. Variance explained by each Principal Component (PC) up to the 20th PC. This type of a plot is also known as a Scree Plot

3.2 Principal Components 1 and 2

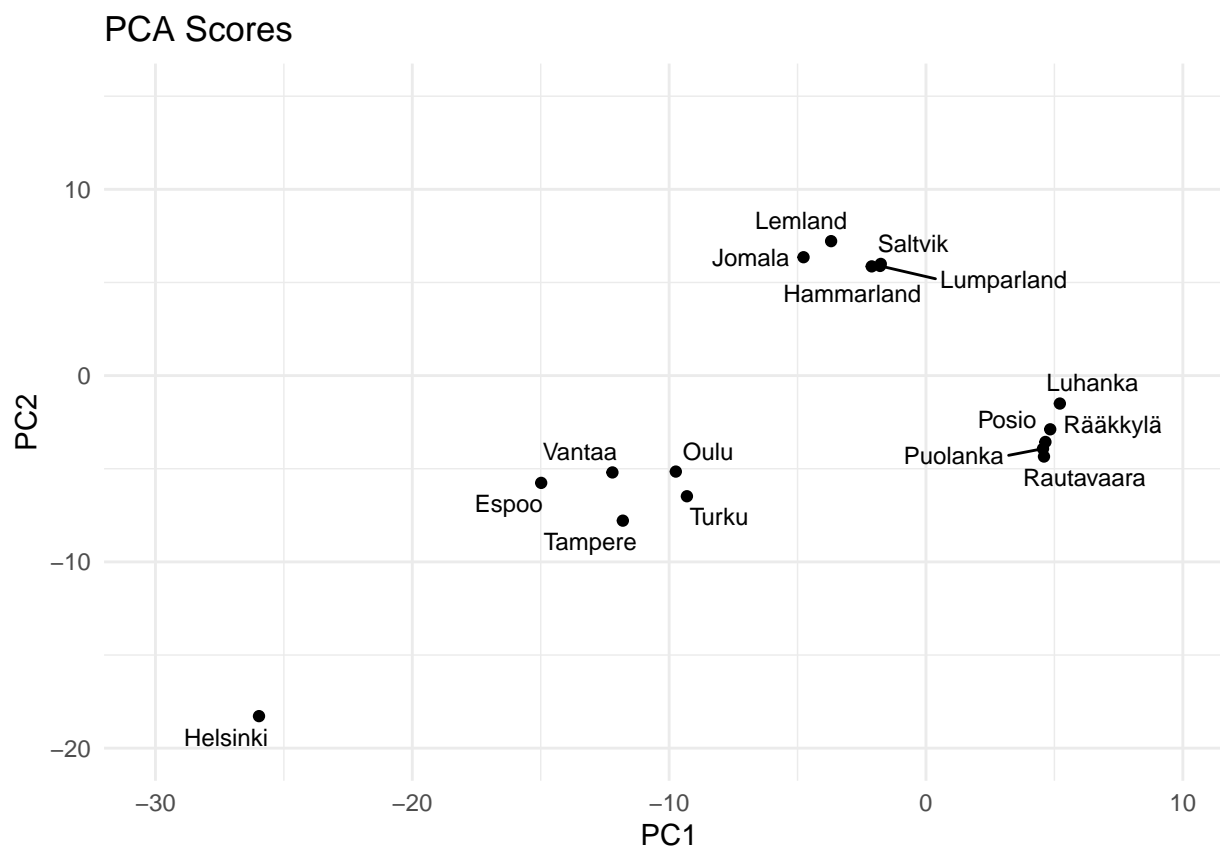


Figure 7. Score plot for all municipalities that have a top 5 absolute loading in any of PC1 and PC2

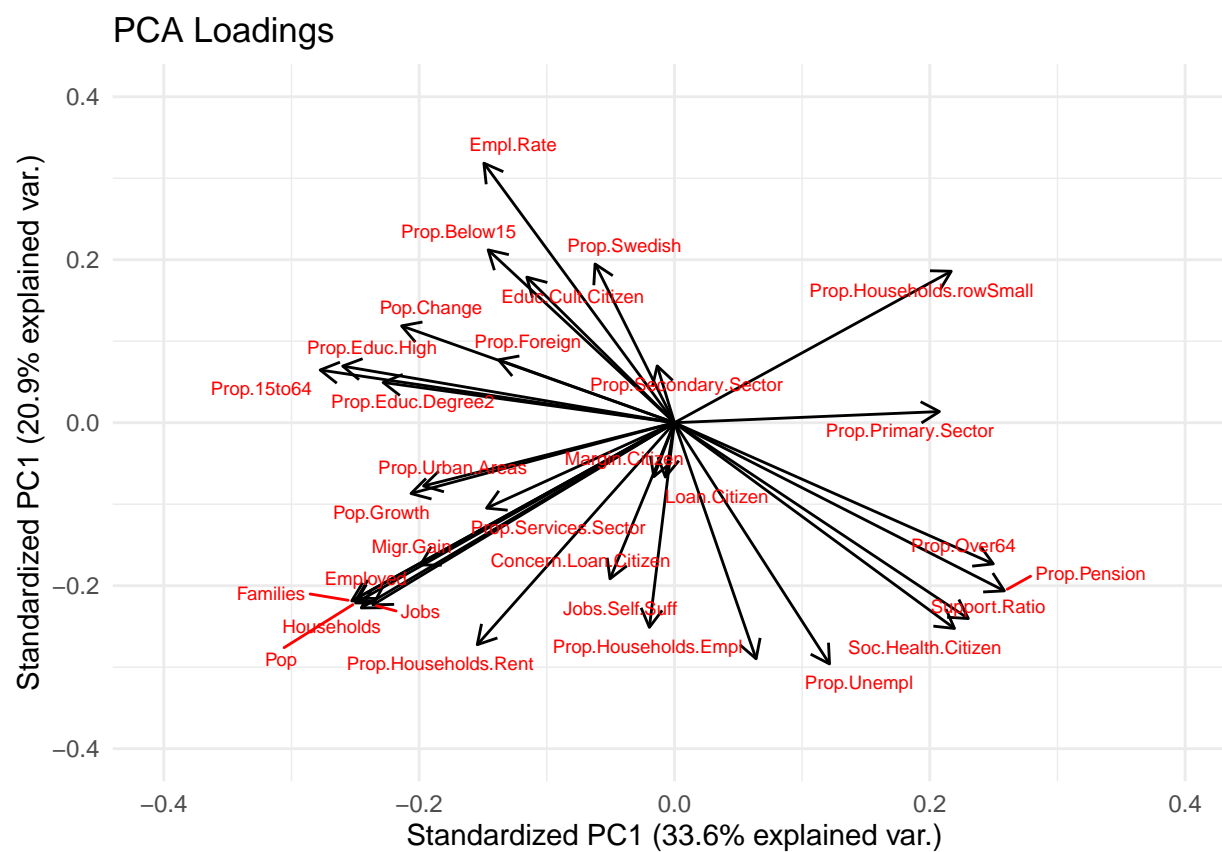


Figure 8. Loading directions for PC1 and PC2

3.3 Principal Components 3 and 4

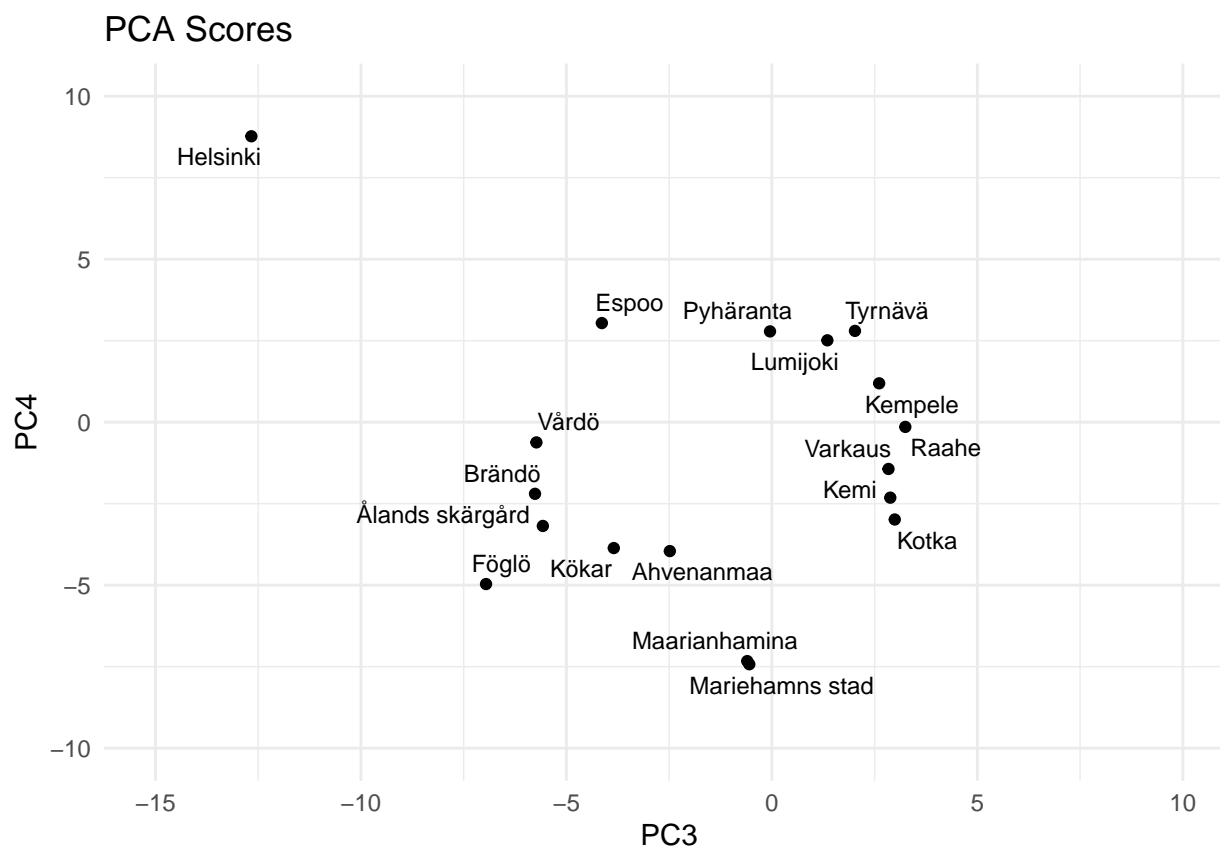


Figure 9. Score plot for all municipalities that have a top 5 absolute loading in any of PC3 and PC4

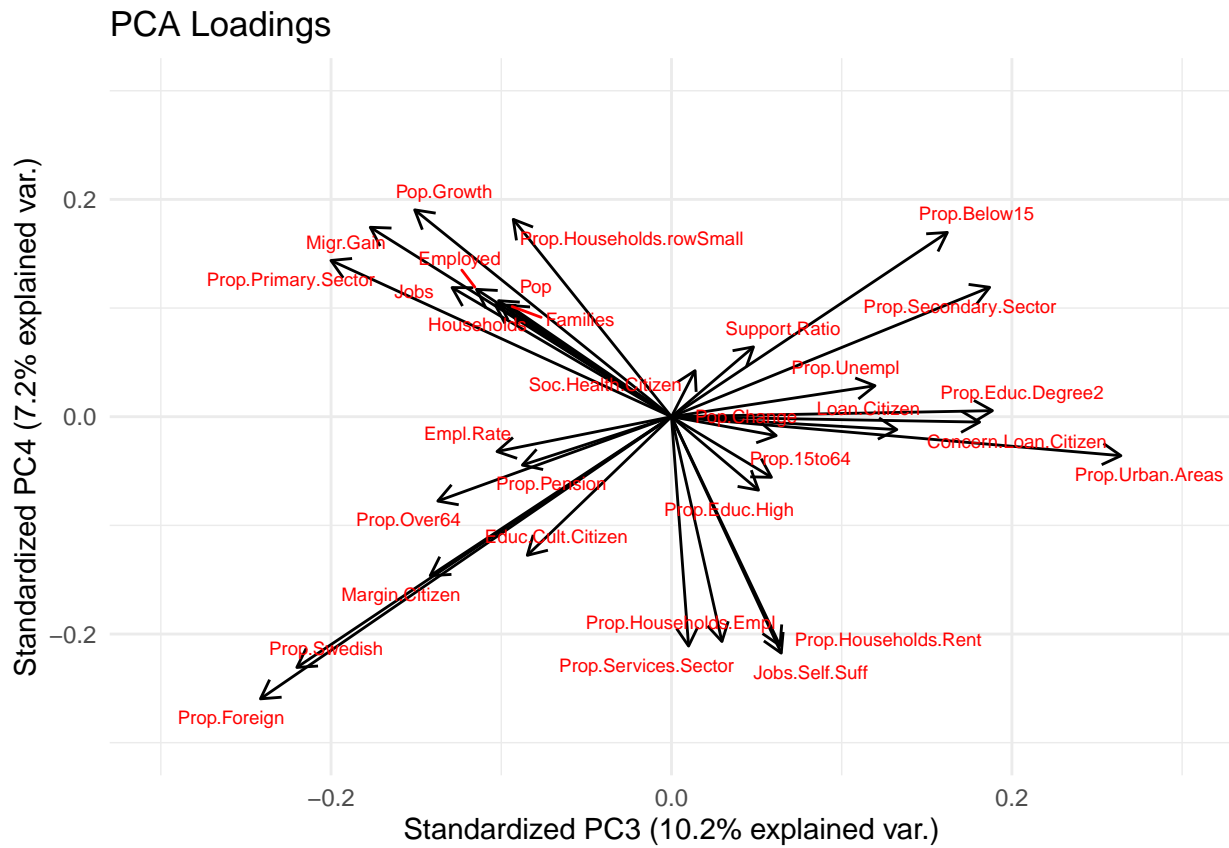


Figure 10. Loading directions for PC3 and PC4

4 Discussion and Conclusions

- PC1 and PC2 do not explain a high percentage of the variance (above 80%). Hence, PCA might not be so well suited
- Using population, employed persons and some other absolute measures might have alienated Helsinki from the other cities and biased the analysis a bit. Perhaps leaving out Helsinki would have been a wise choice
- Something else?

References

Mans Magnusson, Markus Kainu, Janne Huovari, and Leo Lahti. pxweb: R tools for px-web api, 2019.

Statistics Finland. Kuntien avainluvut muuttujina alue 2020, tiedot ja vuosi, 2022. URL https://pxwebapi2.stat.fi/PXWeb/api/v1/fi/Kuntien_avainluvut/2020/kuntien_avainluvut_2020_aikasarja.px. [Data accessed 2022-04-06 16:59:07 using pxweb R package 0.13.1].