

Study of French labour market and inequalities

L. Insolia, J. Kim and G. Yeghikyan

SNS

— *Midterm results* —

March 14, 2018

Objectives

- Structure of French labour market
- Inequalities (in terms of salary):
 - ages
 - gender
 - job categories
 - spatial distribution
- Firms' distribution
- Exploratory analyses

Methodology

INSEE data

- Population: age, sex and cohabitation mode
- Salary: job categories, age and sex (mean net salary per hour in €)
- Firms: number of firms for each size
- Geography: GPS location

for different geographical levels (communes, departments, towns) in 2014

GitHub repo: <https://github.com/LucaIns/TSL>

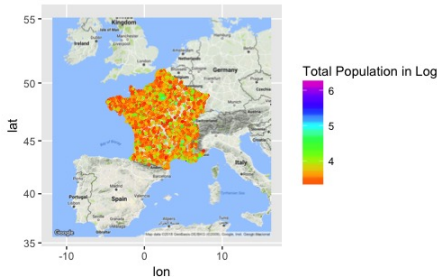
What has been done so far . . .

Pre-processing phase

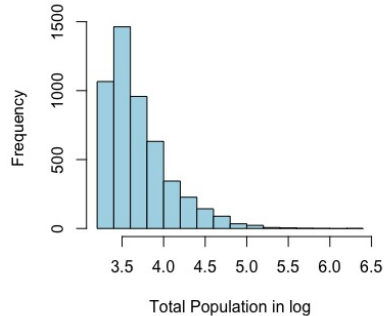
- Population: restructured the dataset and created new features
- Firms: categorized firms' sizes into 4 categories
- Geography: retrieved the missing data using Google API

Distribution of population per town

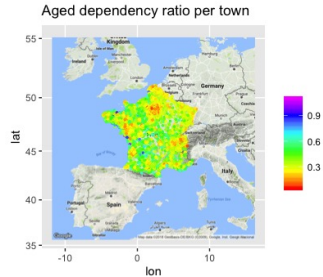
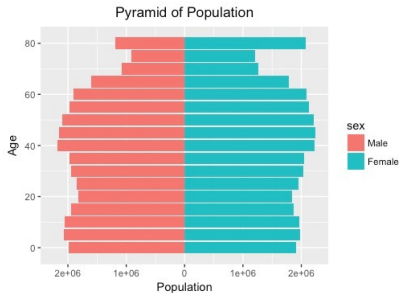
Distribution of Population for each town



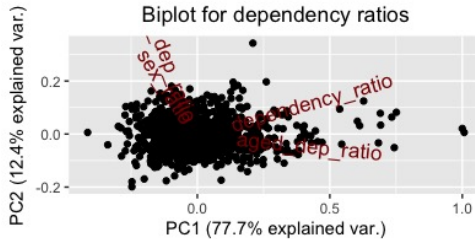
Total population per town



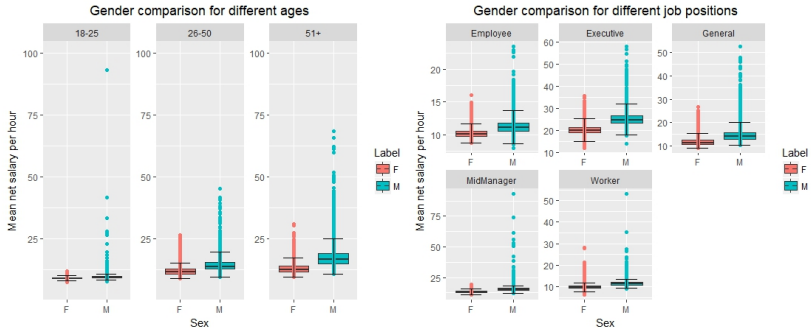
Population demographics



PCA

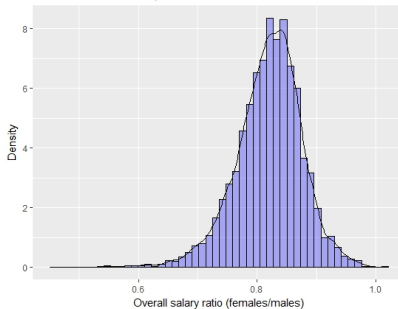


Inequality of salary

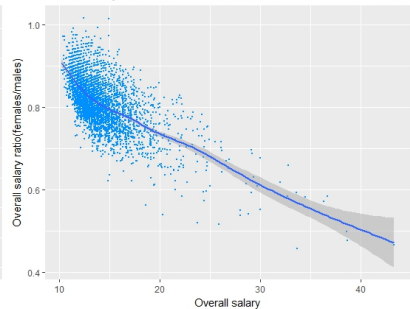


Inequality of salary

Overall salary ratio between females and males

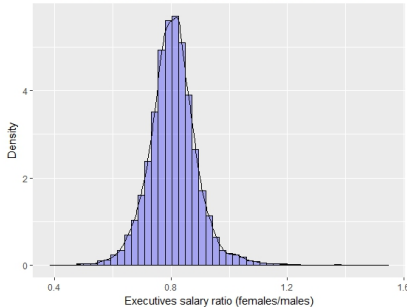


Overall salary ratio between females and males vs. overall salary



Inequality of salary

Executives salary ratio between females and males



Executives salary ratio between females and males
vs. overall executives salary

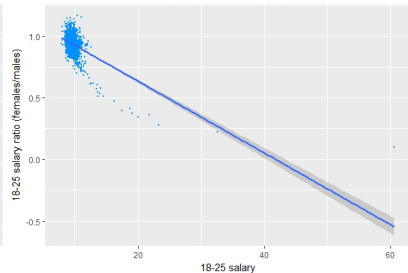


Inequality of salary

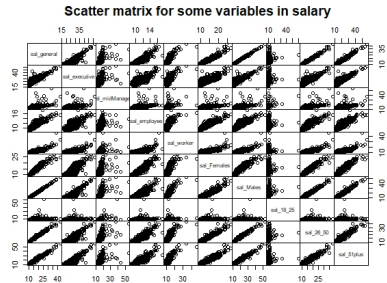
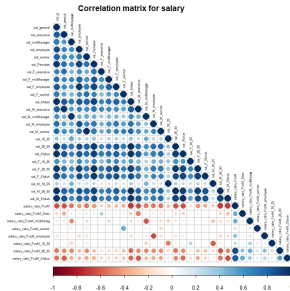
18-25 salary ratio between females and males
vs. overall salary



18-25 salary ratio between females and males
vs. overall 18-25 salary

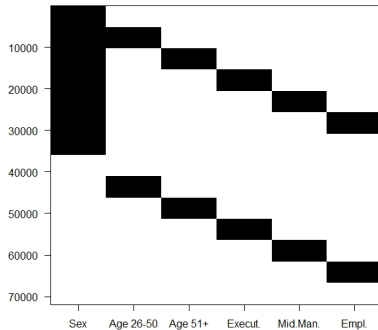


Bivariate relations

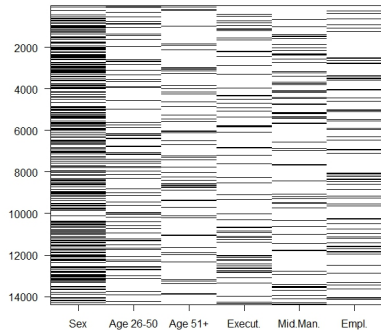


ANOVA

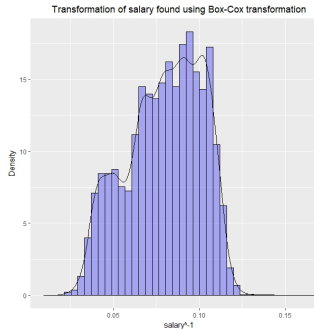
Factors for ANOVA



Factors for ANOVA



ANOVA



```
Call:
lm(formula = sal_y ~ sex + age + job + sex:age + sex:job)

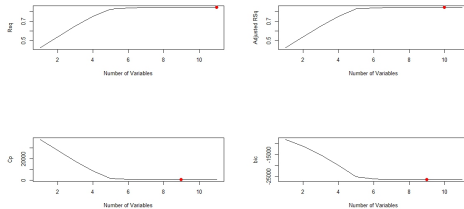
Residuals:
    Min       1Q   Median       3Q      Max
-0.070217 -0.004526  0.000711  0.005522  0.057875

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.1060590  0.0001889   561.599 < 2e-16 ***
sex          -0.0106241  0.0002657   -39.989 < 2e-16 ***
age1         -0.0216544  0.0003289   -65.833 < 2e-16 ***
age2         -0.0278782  0.0003254   -85.667 < 2e-16 ***
job1         -0.0560085  0.0003239  -172.923 < 2e-16 ***
job2         -0.0303178  0.0003275   -92.566 < 2e-16 ***
job3         -0.0085416  0.0003254   -26.247 < 2e-16 ***
sex:age1     -0.0033048  0.0004648    -7.110 1.22e-12 ***
sex:age2     -0.0084595  0.0004613   -18.339 < 2e-16 ***
sex:job1      0.0005974  0.0004607     1.297  0.195
sex:job2      0.0002832  0.0004626     0.612  0.540
sex:job3      0.0026619  0.0004639     5.738 9.75e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

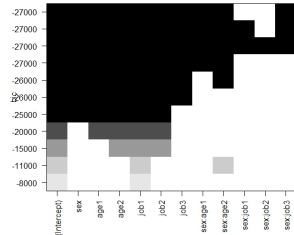
Residual standard error: 0.008563 on 14369 degrees of freedom
Multiple R-squared:  0.8432,    Adjusted R-squared:  0.8431
F-statistic: 7026 on 11 and 14369 DF, p-value: < 2.2e-16
```

ANOVA BSS

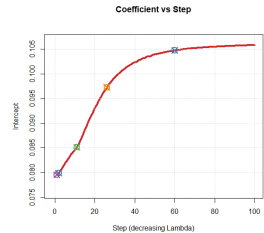
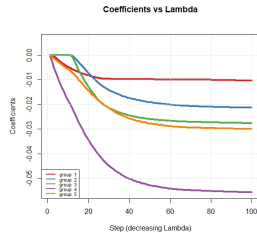
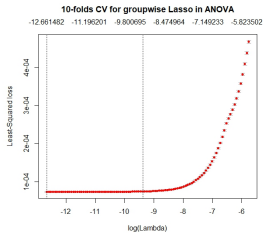
Best subset selection for ANOVA



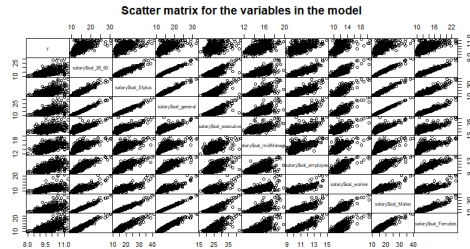
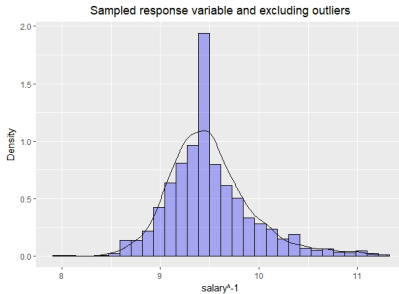
Best subset selection for ANOVA



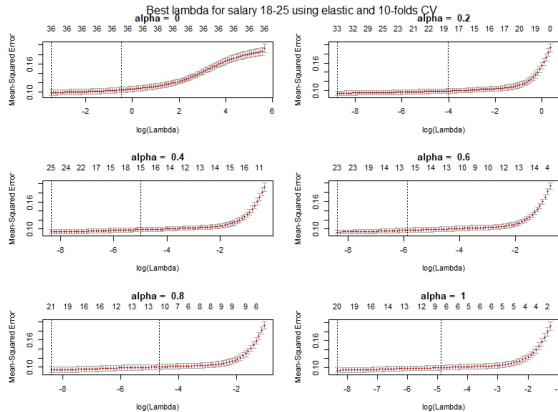
ANOVA GW-Lasso



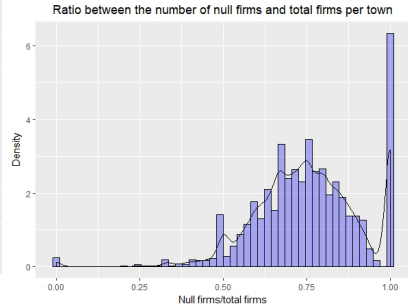
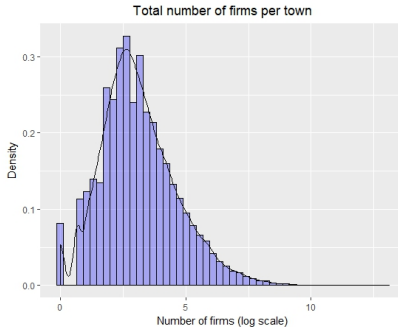
Prediction for young people



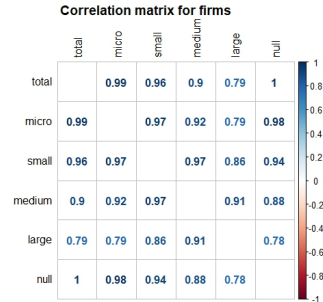
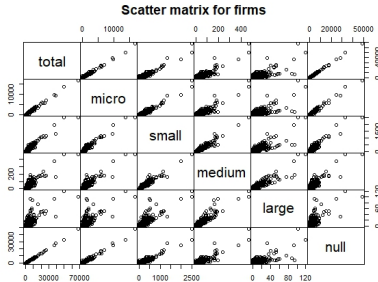
Elastic net and 10-folds CV



Distribution of firms per town



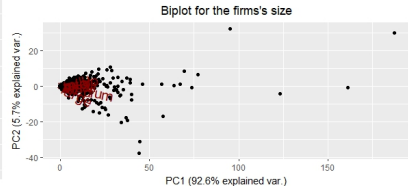
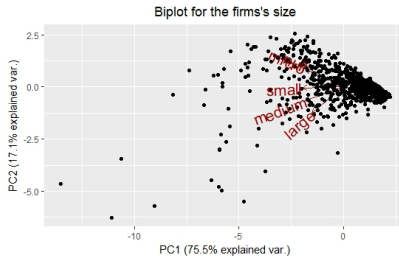
Bivariate relations



Excluding Paris

PCA

Using original data scaled (not logs)
Most typical vs. Excluding just Paris



Issues

- Unique code for salary data 1/7 of the total
- Loss of information when combining the separated datasets
- Missing additional information
- French DOM-TOM regions
- Outliers and spatial correlation

Future works

- Create meaningful indicators
- Clustering techniques to identify geographical clusters
- Verification and improvement of the obtained results
- Compare the methodologies used with robust ones
- Find complementary datasets

– *Thank you* –