



INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL

» Mismatch between jobs and skills in the EU

JOCLAD 2017 «

DMSI-ME / Joao S. Lopes

■ ■ ■ ■ (21-04-2017)



Motivation



- “Skills development are essential in the emerging new economy and fast-changing labour market”¹
- “Qualification and skill mismatches entail significant economic and social costs for individuals and firms”¹
- Skills mismatch (i.e. over-qualification, under-qualification) remains at 45% (CEDEFOP, 2015)²
- EU Guidelines (2015) call for enhancing labour supply, skills and competences³

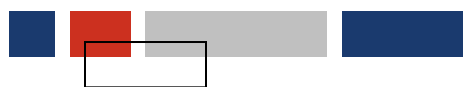


Motivation



Create framework that:

1. combines **Official Statistics** with **Big Data**
2. estimates **Labour Market Attractiveness** and its association with **Skills Mismatch**, Labour Market Mobility and Emigration
3. is aimed at **policy makers** and both **jobseekers** and **job providers**

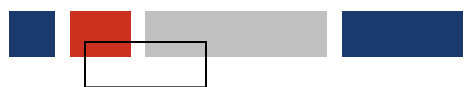


Data: LMkt Attractiveness



- “reg_dem” – demographic statistics
- “earn” – earning structure
- “educ_uoe_fin” – public expenditure on education
- “ilc” – income and life conditions
- “employ” – employment statistics
- “nama10” – annual national accounts
- “educ_part” – participation in education

7 datasets, 17 main variables



Data: LMkt Attractiveness

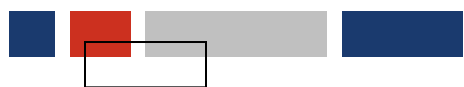


- “reg_dem” by **age** (**NUTS2**)
- “earn” by **occupation** and **economic activity**
- “educ_uoe_fin”
- “ilc” (**NUTS2**)
- “employ” by **age**, **education level**, **economic activity** (**NUTS2**)
- “nama10” (**NUTS2**)
- “educ_part” (**NUTS2**)

7 datasets, **17** main variables, **76** variables

subjects: **NUTS0 = 28**; NUTS1= 98; NUTS2 = 276.





Data: Skills mismatch



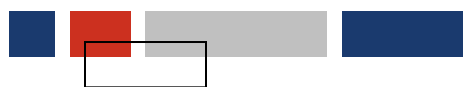
- “EURES” **scrapped data** on jobseekers’ CVs
- “EURES” **scrapped data** on job vacations

2 datasets, **1** main variable

... but

cleaning and **structuring** requires considerable expertise

normalization requires detailed demographical information

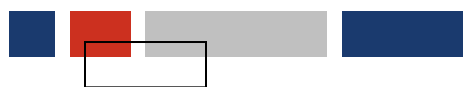


Data: Skills mismatch



- “educ_uoe_grad02” – education statistics
- “jvs_a_nace2” – job vacancy statistics

2 datasets, 1 main variable



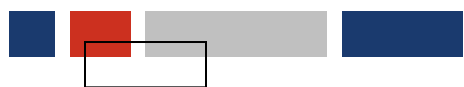
Data: Skills mismatch



- “educ_uoe_grad02” by **education field**
- “jvs_a_nace2” by **occupation** and **economic activity (NUTS2)**

2 datasets, **1** main variable, **1** variable

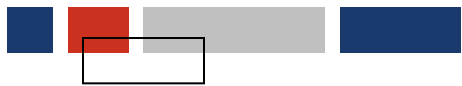
subjects: **NUTS0 = 8**; NUTS1= 14; NUTS2 = 47



Methods



- Network Analysis
- Partition-around-medoids (PAM)⁴
- Over-representation analysis (ORA)
- Multinomial regression
- Multivariate regression
- Weighted correlation network analysis (WCNA)⁵



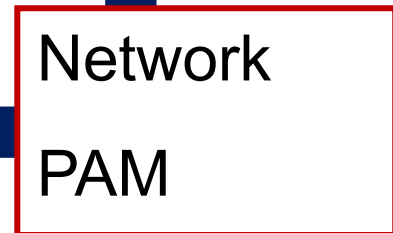
Methods



Labour Market Attractiveness (by NUTS 0-2)



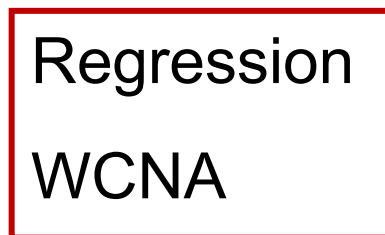
distance between NUTS



NUTS groups

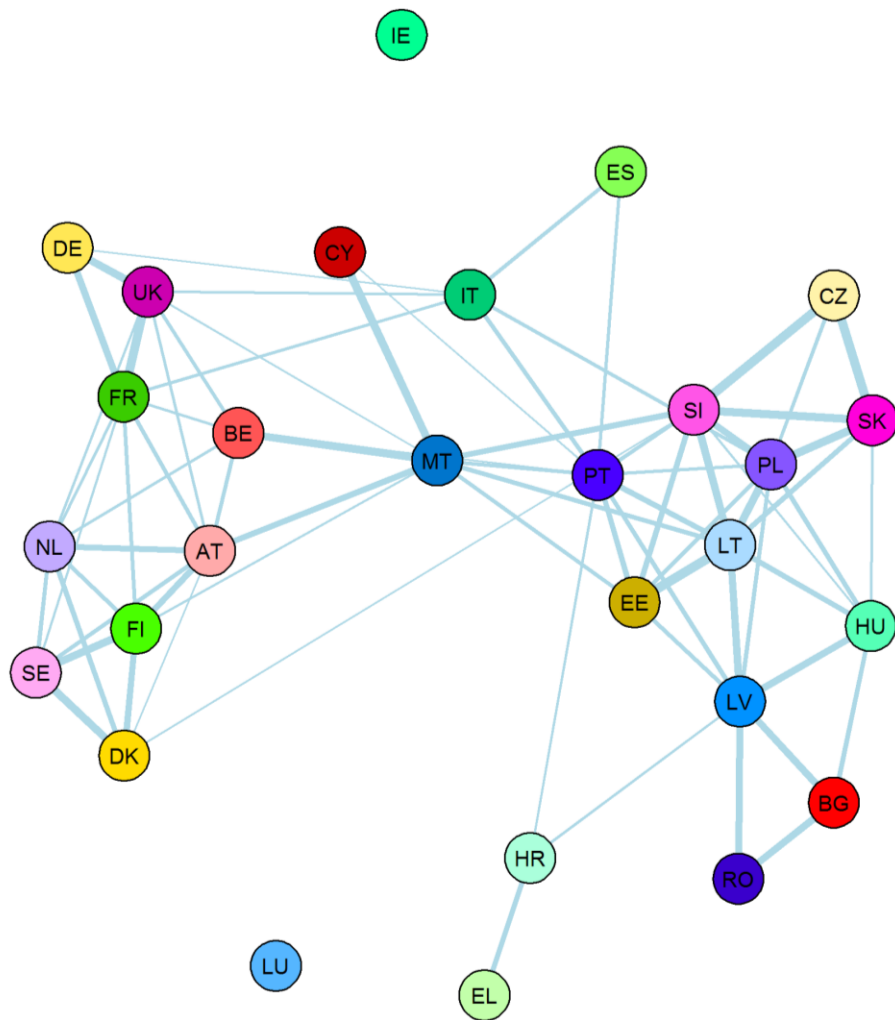


Skills mismatch (by NUTS0-2)

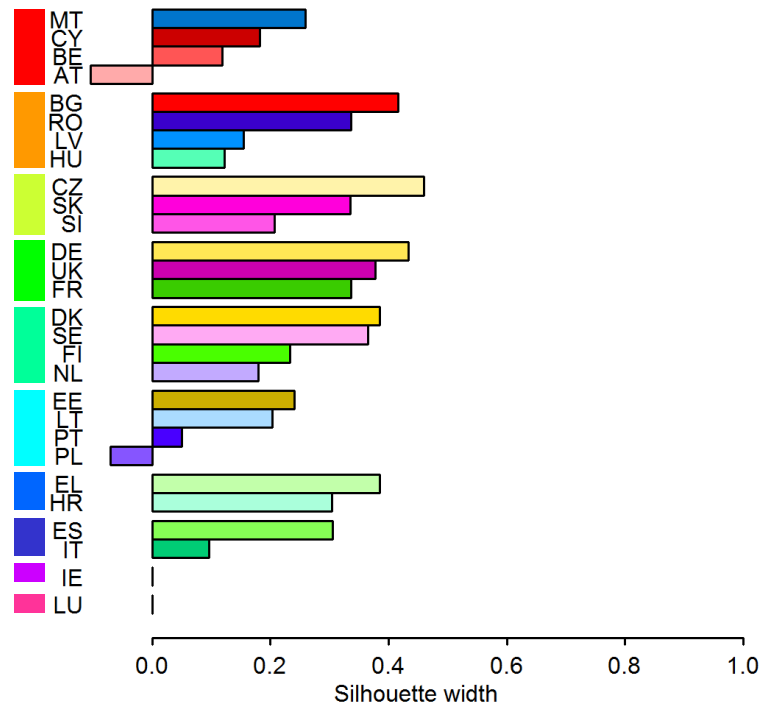




Results: LMkt Attractiveness



Labour market attractiveness





Results: LMkt Attractiveness



“MT-CY-BE-AT”

emp_15-24_ED5-8

expend_ED5-8

“BG-RO-LV-HU”

ARPR_socexcl

emp_Y15-24_NaceA

earn_OC[1-9]

“CZ-SK-SI”

emp_FT

pop_Y25-64

ARPR_socexcl

“DE-UK-FR”

emp_PT

disp_income

pop_Total

“DK-SE-FI-NL”

earn_OC[1-9]

expend_ED5-8

ARPR_socexcl

“EE-IT-PT-PL”

earn_OC[1-5,7-9]

“EL-HR”

GVAgr

“ES-IT”

emp_Y15-24_Nace[A,J-L]

expend_ED5-8

“IE”

GDP

GVAgr

pop_Y0-14

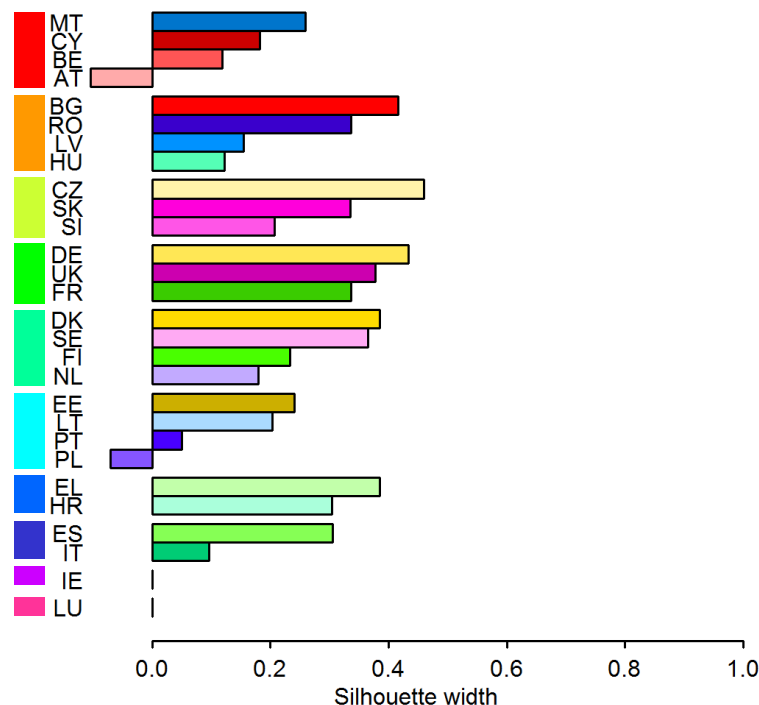
“LU”

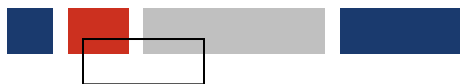
emp_Y25-64_ED5-8

GDP

low_work

Labour market attractiveness

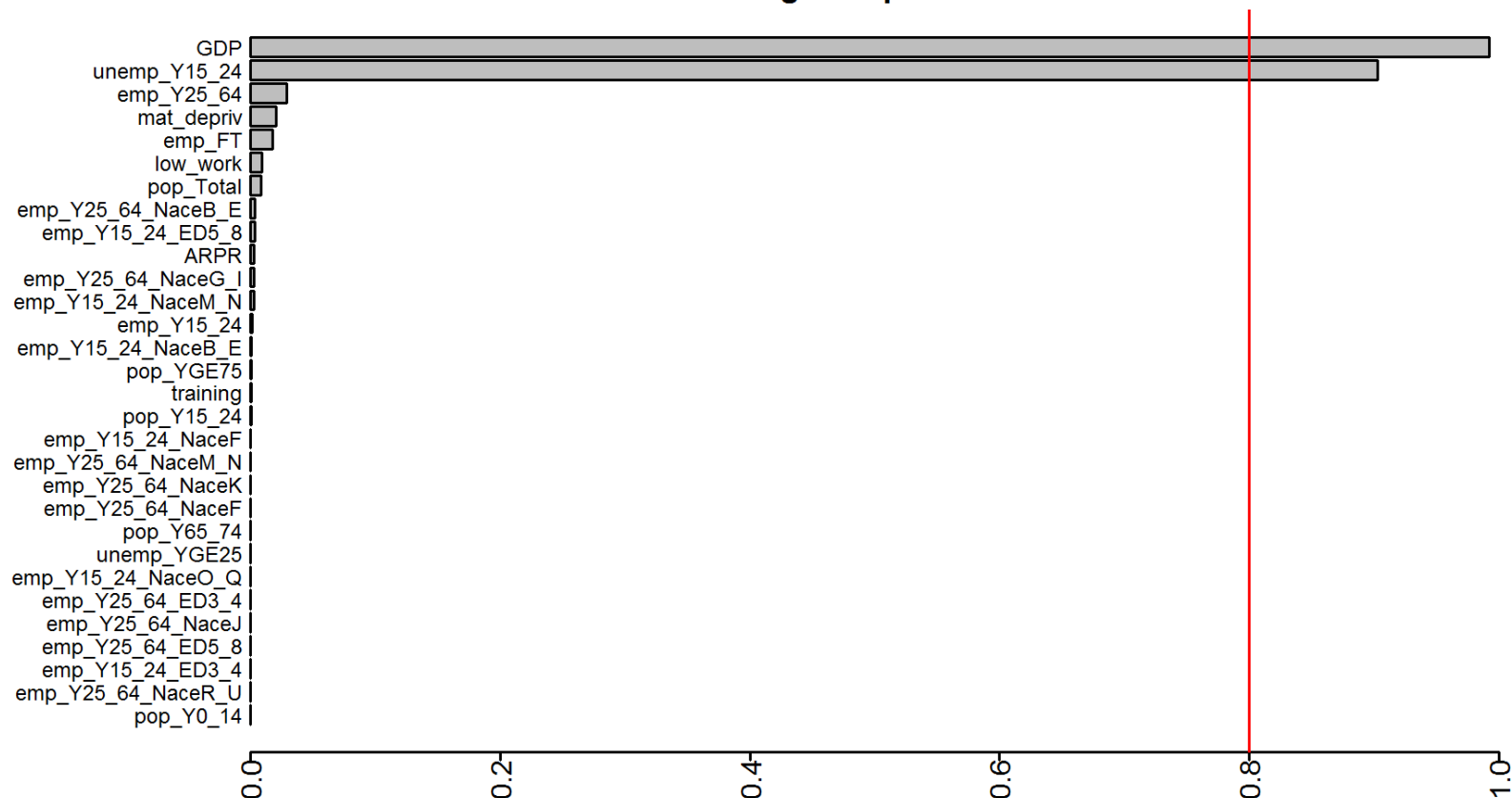




Results: LMkt Attractiveness



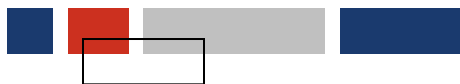
Model-averaged importance of terms



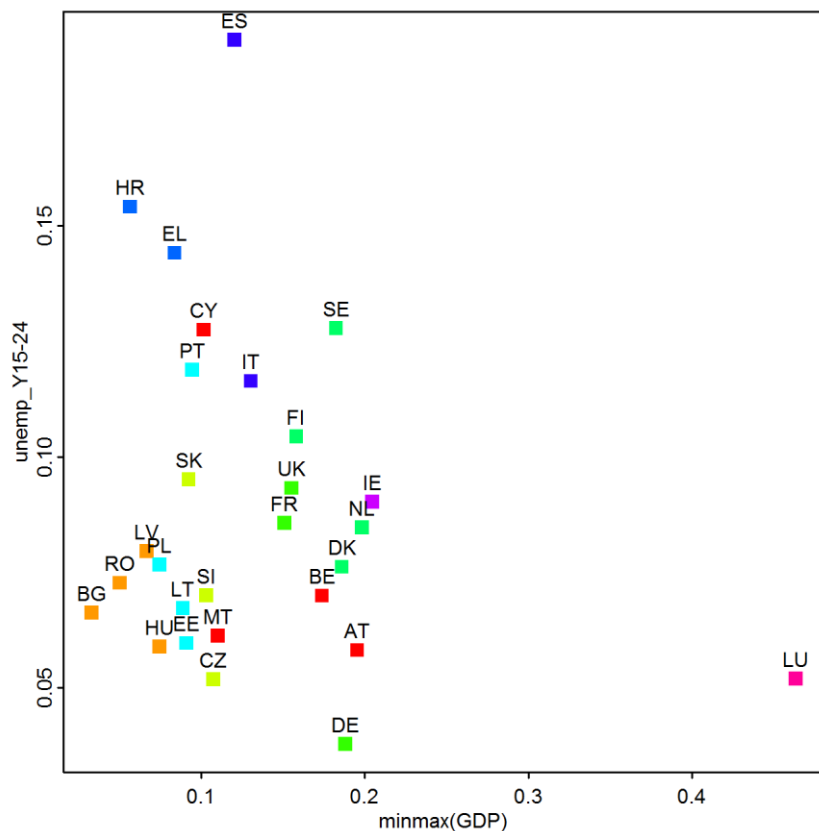
EU_groups ~ 1 + unemp_Y15-24 + GDP

R_{McFadden} = 0.82, R_{count} = 0.82

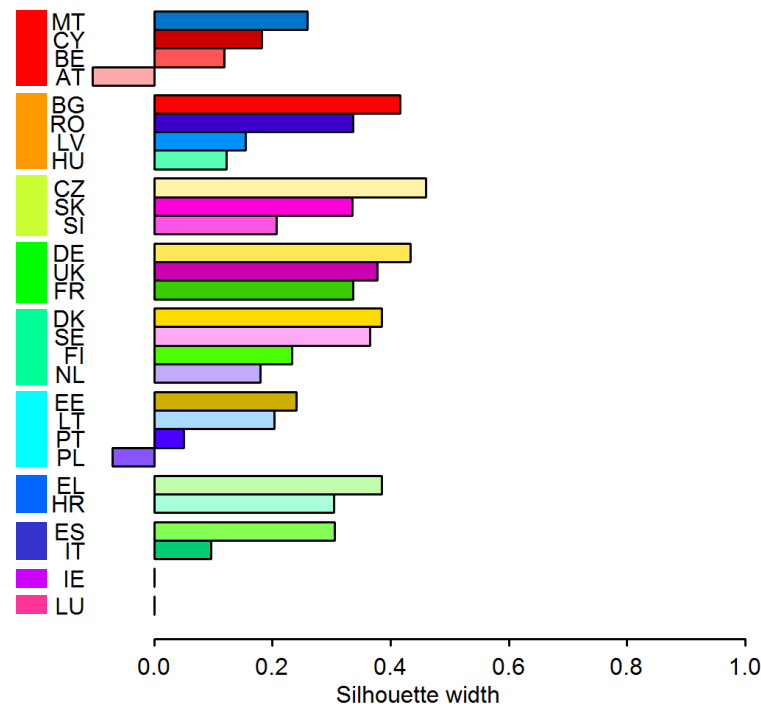




Results: LMkt Attractiveness



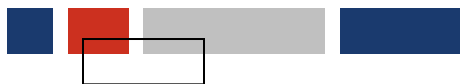
Labour market attractiveness



EU_groups ~ 1 + unemp_Y15-24 + GDP

R_{McFadden} = 0.82, R_{count} = 0.82

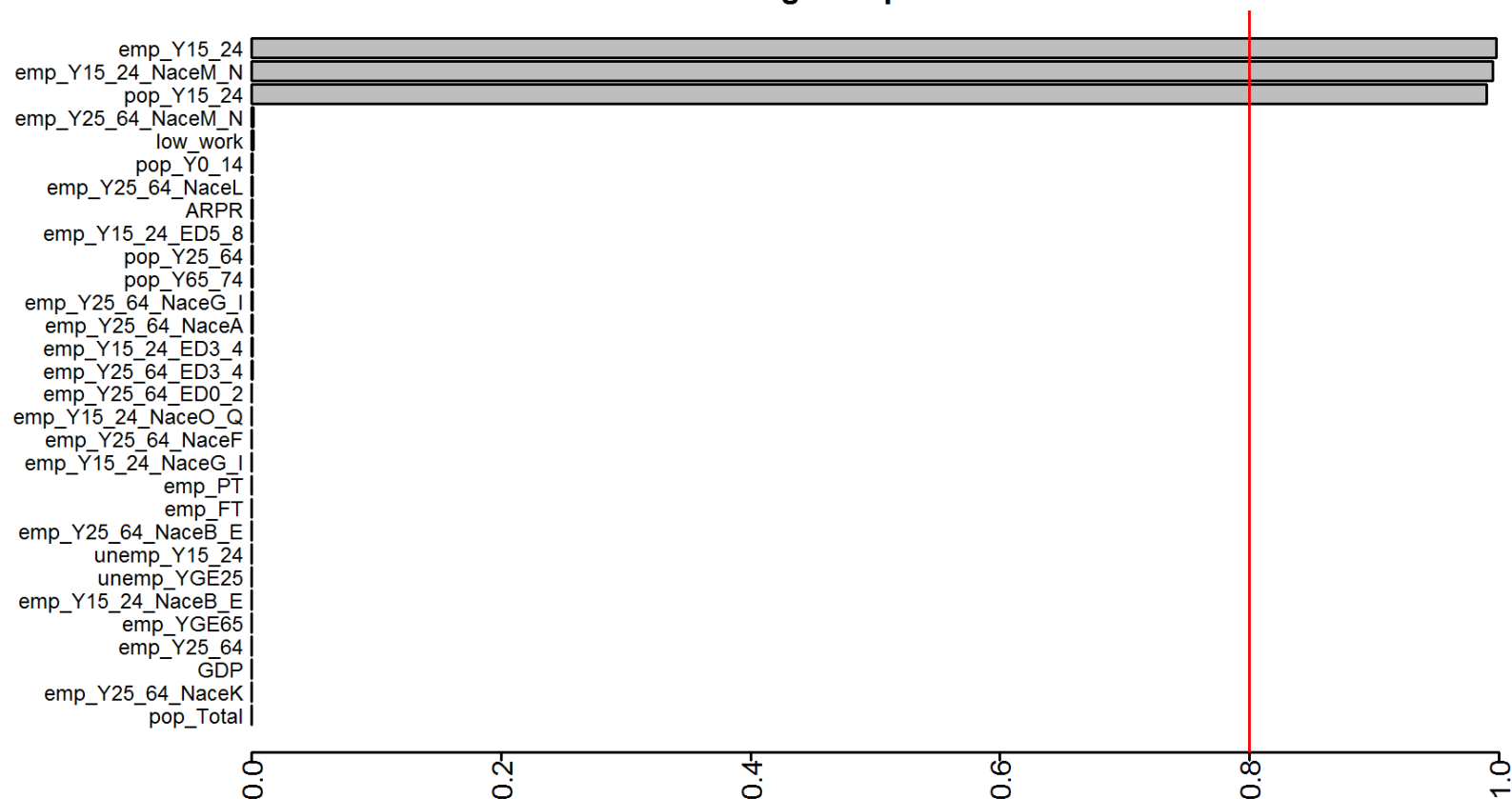




Results: Skills mismatch

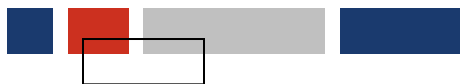


Model-averaged importance of terms

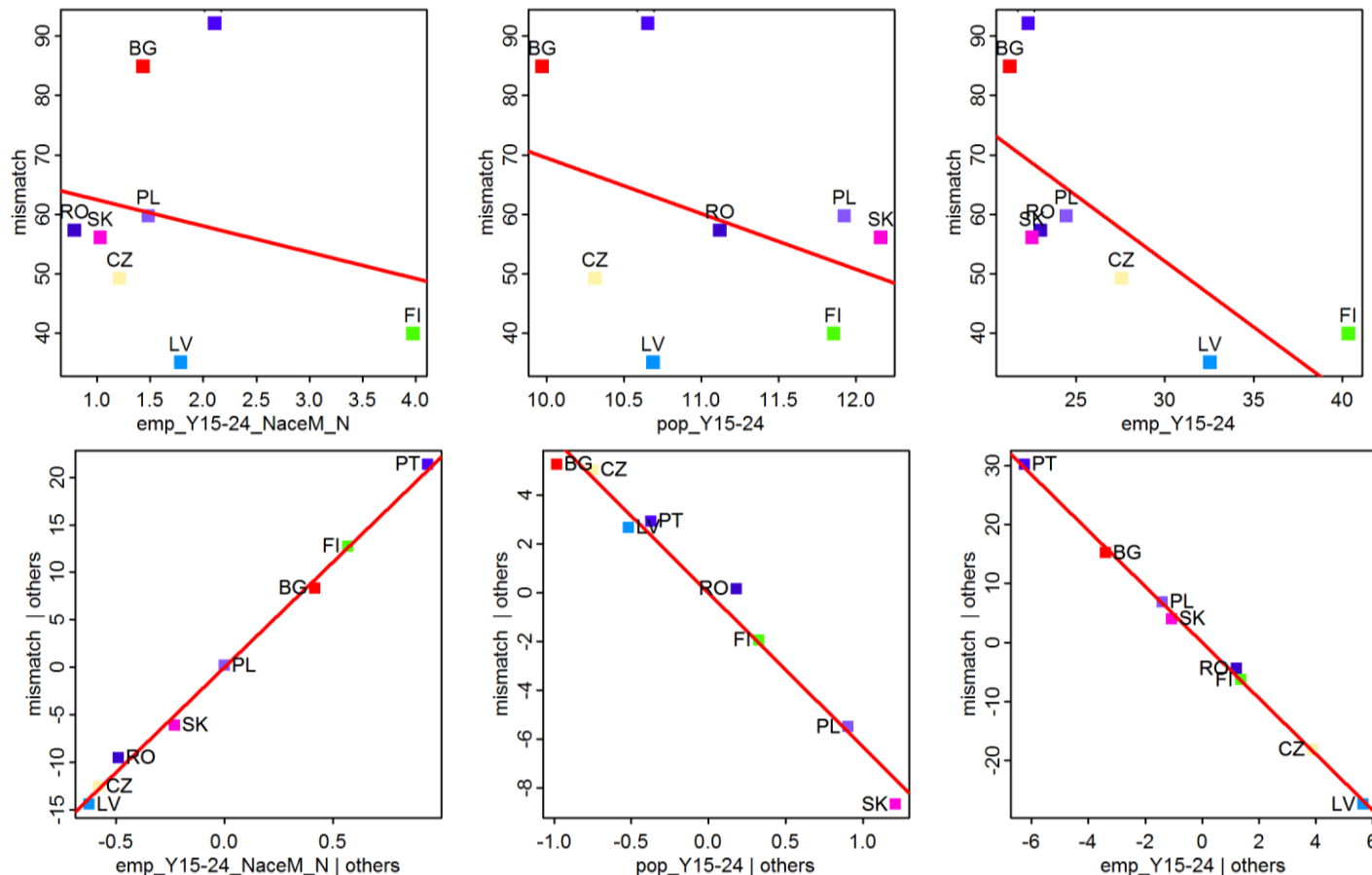


$$\text{mismatch} = 217.8 + 22.2 \cdot \text{emp_Y15-24_NaceM_N} - (6.3 \cdot \text{pop_Y15-24} + 4.7 \cdot \text{emp_Y15-24})$$

$$R^2 = 1.00$$



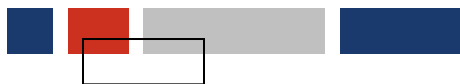
Results: Skills mismatch



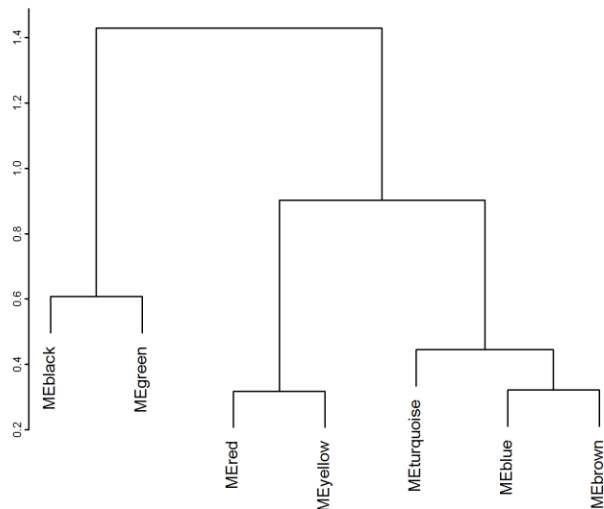
$$\text{mismatch} = 217.8 + 22.2 \cdot \text{emp_Y15-24_NaceM_N} - (6.3 \cdot \text{pop_Y15-24} + 4.7 \cdot \text{emp_Y15-24})$$

$$R^2 = 1.00$$





Results: WCNA



MEblack

unemp_Y[15-24,25-64]

MEgreen

ilc_ARPR

ilc_ARPR_socexcl

pop_YG75

MEred

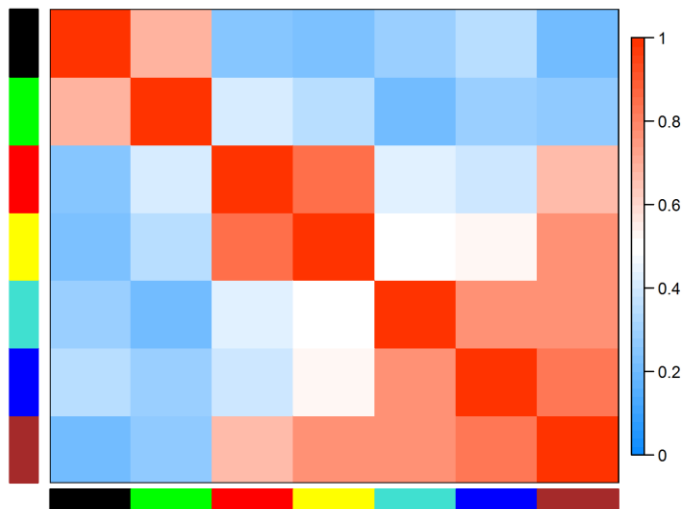
emp_Y[15-24,25-64]_NaceB-E

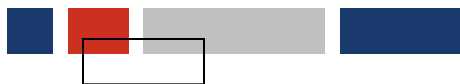
MEyellow

ilc_low_work

emp_Y25-64_ED[0-2,3-4]

emp_Y25-64_NaceF





Results: WCNA



MEturquoise

ilc_mat_depriv

ilc_rooms_pp

earn_OC[1-9]

emp_[FT,PT]

emp_Y15-24_ED0-2

emp_Y15-24_Nace[A,O-Q,R-U]

emp_Y25-64_ED5-8

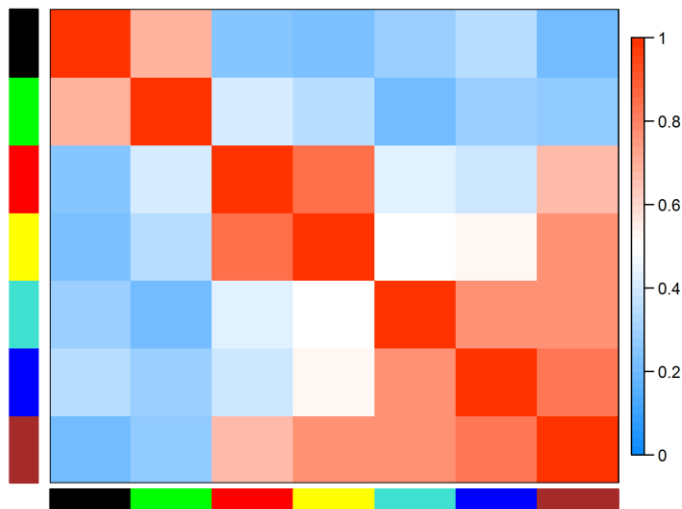
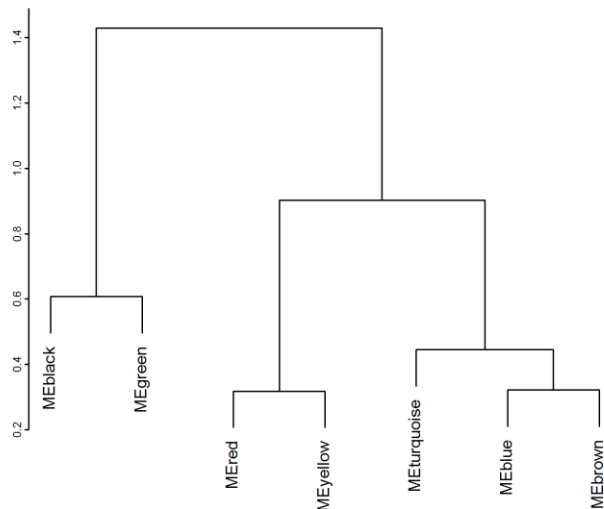
emp_Y25-64_Nace[A,G-I,J,K,M_N,O-Q,R-U]

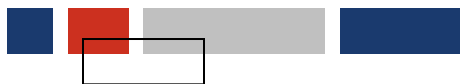
na_disp_income

na_GDP

pop_Y[0-14,25-64]

training





Results: WCNA



MEblue

emp_Y15-24_ED5-8

emp_Y15-24_Nace[G-I,J,M_N]

emp_YGE65

pop_Y15-24

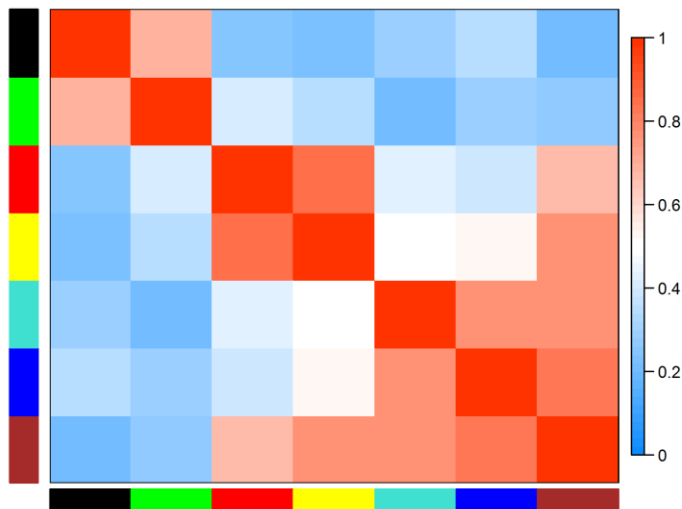
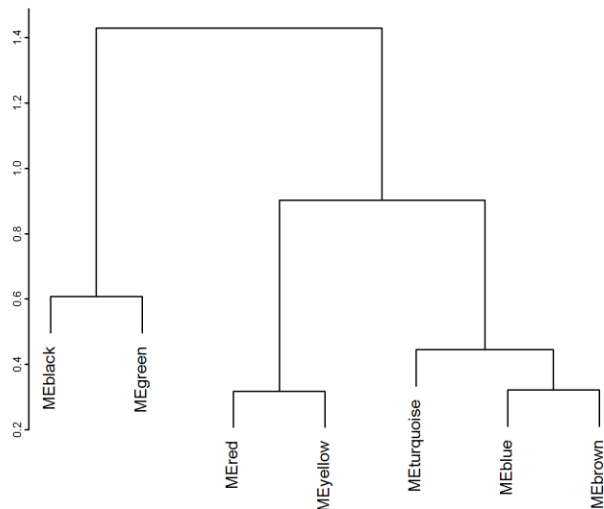
MEbrown

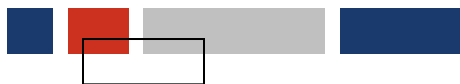
emp_Y[15-24,25-64]

emp_Y15-24_ED3-4

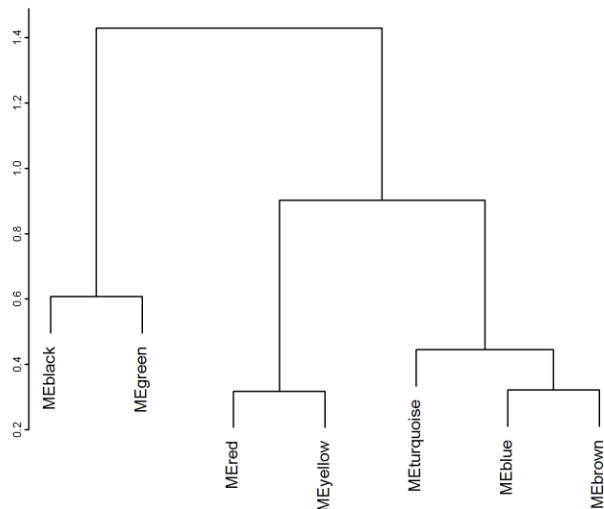
emp_Y15-24_NaceF

emp_Y25-64_NaceL

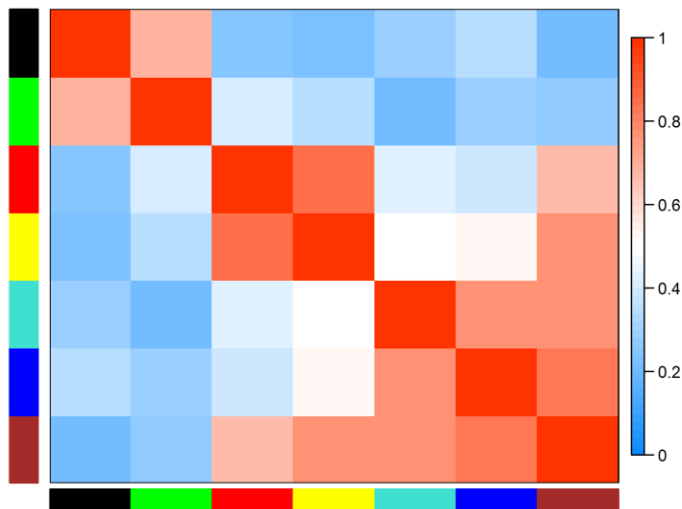


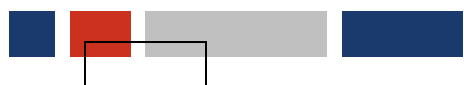


Results: WCNA



labels	description	mismatch	EU_groups
MEblack	Unemployment		
MEgreen	Poverty		0.63
MEred	Industry sector		
MEyellow	emp_Y25-64_ED0-4	-0.50	
MEturquoise	Earnings		0.80
MEblue	emp_Y15-24_ED5-8		0.63
MEbrown	emp_Y15-24_ED3-4	-0.86	0.57





Conclusions



- **LMkt Attract** is able to form consistent clusters at NUTS0
- **Youth unemployment** and **GDP** can separate well clusters
- **LMkt Attract** can be separate in different modules: **Unemployment**, **Poverty**, **Industry**, **emp_Y25-64_ED0-4**, **Earnings**, **emp_Y15-24_ED3-4** and **emp_Y15-24_ED5-8**
- **Skills Mismatch** is associated with **population Y15-24**, negative association with **Pop_{prop}**, **Emp_{prop}** and **NaceM_N_{prop}**
- **Skills Mismatch** is strongly associated to module **emp_15-24_ED3-4**



Acknowledge



Team:

Sónia Quaresma

João Lopes

Marco Moura

Institution:



Data:



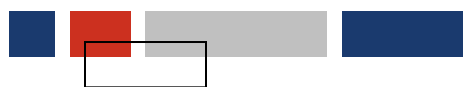
INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL



Thank you!

JOCLAD 2017





Bibliography



1. https://ec.europa.eu/commission/publications/skills-education-and-lifelong-learning-european-pillar-social-rights_en
2. CEDEFOP (2015) “Skills, qualifications and jobs in the EU: the making of a perfect match? “
3. Council Decision (EU) 2015/1848 of 5 October 2015
4. Reynolds et al. (1992) “Clustering rules: A comparison of partitioning and hierarchical clustering algorithms” J Math. Model. Algorithms
5. Langfelder and Horvath (2008) “WGCNA: an R package for weighted correlation network analysis” BMC Bioinformatics



R libraries

car - Companion to Applied Regression

caret - Classification and Regression Training

cluster - Finding Groups in Data: Cluster Analysis Extended

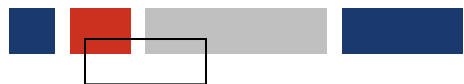
glmulti - Model selection and multimodel inference made easy

MASS - Support Functions and Datasets for MASS

nnet - Feed-Forward Neural Networks and Multinomial Log-Linear Models

sna - Tools for Social Network Analysis

WGCNA - Weighted Correlation Network Analysis



Metadata: ISCED 11



label	description
ED0-2	Less than primary, primary and lower secondary education (levels 0-2)
ED3_4	Upper secondary and post-secondary non-tertiary education (levels 3 and 4)
ED5-8	Tertiary education (levels 5-8)



Metadata: ISCED-F 13



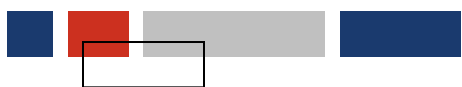
label	description
F00	Generic programmes and qualifications
F01	Education
F02	Arts and humanities
F03	Social sciences, journalism and information
F04	Business, administration and law
F05	Natural sciences, mathematics and statistics
F06	Information and Communication Technologies
F07	Engineering, manufacturing and construction
F08	Agriculture, forestry, fisheries and veterinary
F09	Health and welfare
F10	Services



Metadata: ISCO-08



label	description
OC1-5	Non manual workers
OC1	Managers
OC2	Professionals
OC3	Technicians and associate professionals
OC4	Clerical support workers
OC5	Service and sales workers
OC6-8	Skilled manual workers
OC6	Skilled agricultural, forestry and fishery workers
OC7	Craft and related trades workers
OC8	Plant and machine operators and assemblers
OC9	Elementary occupations
OC0	Armed forces occupations



Metadata: NACE Rev. 2



label	description
A	Agriculture, forestry and fishing
B-E	Industry (except construction)
B-F	Industry and construction
B-N	Business economy
F	Construction
G-I	Wholesale and retail trade, transport, accomodation and food service activities
G-N	Services of the business economy
J	Information and communication
K	Financial and insurance activities
L	Real estate activities
M_N	Professional, scientific and technical activities; administrative and support service activities
O-Q	Public administration, defence, education, human health and social work activities
P-S	Education; human health and social work activities; arts, entertainment and recreation; other service activities
R-U	Arts, entertainment and recreation; other service activities; activities of household and extra-territorial organizations and bodies



Data: LMkt attractiveness



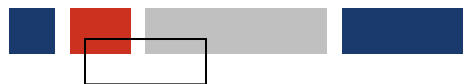
dataset	description	year	NUTS	units
demo_r_d2jan	Population	2015	NUTS 2	NR
earn_ses_hourly	Structure of earnings: hourly earnings	2014	NUTS 0	MN_PPS
educ_uoe_fine06	Total public expenditure on education	2013	NUTS 0	PC_GDP
ilc_li41	At-risk-of-poverty rate	2014	NUTS 2	PC_POP
ilc_lvhl21	People living in households with very low work intensity	2014	NUTS 2	PC_Y_LT60
ilc_lvho04n	Average number of rooms	2014	NUTS 2	AVG
ilc_mddd21	Severe material deprivation rate	2014	NUTS 2	PC_POP
ilc_peps11	People at risk of poverty or social exclusion	2014	NUTS 2	PC_POP
lfst_r_lfe2eedu	Employment by educational attainment level (ISCED 11)	2014	NUTS 2	THS
lfst_r_lfe2eftpt	Employment by full-time/part-time	2014	NUTS 2	THS
lfst_r_lfe2emp	Employment	2014	NUTS 2	THS
lfst_r_lfe2en2	Employment by economic activity (NACE Rev. 2)	2014	NUTS 2	THS
lfst_r_lfu3pers	Unemployment	2014	NUTS 2	THS
nama_10r_2gdp	Gross domestic product	2014	NUTS 2	PPS_HAB
nama_10r_2gvagr	Real growth rate of regional gross value added	2014	NUTS 2	PCH_PRE
nama_10r_2hhinc	Income of households	2013	NUTS 2	PPCS_HAB
trng_lfse_04	Participation rate in education and training (last 4 weeks)	2014	NUTS 2	PC



Data: LMkt attractiveness



variable	description	NUTS	units
ilc_ARPR	At-risk-of-poverty	NUTS 2	PC_POP
ilc_ARPR_socexcl	At-risk-of-poverty or social exclusion	NUTS 2	PC_POP
ilc_low_work	Very low work intensity	NUTS 2	PC_YLT60
ilc_mat_depriv	Severe material deprivation	NUTS 2	PC_POP
ilc_rooms_pp	Number of rooms per person	NUTS 2	AVG
earn_OC[titles]_Nace[sector]	Earning by ISCO-08 title and NACE Ver. 2 sector	NUTS 0	MN_PPS
emp_[contract]	Employment by work contract	NUTS 2	PC_YGE15
emp_Y[age]	Employment by age	NUTS 2	PC_POP[age]
emp_Y[age]_ED[level]	Employment by age and ISCED 11 level	NUTS 2	PC_POP[age]
emp_Y[age]_Nace[sector]	Employment by age and NACE Ver. 2 sector	NUTS 2	PC_POP[age]
unemp_Y[age]	Unemployment by age	NUTS 2	PC_POP[age]
expend_ED5-8	Public expenditure on education	NUTS 0	PC_GDP
na_disp_income	Disposable income	NUTS 2	PPCS_HAB
na_GDP	Gross Domestic Product	NUTS 2	PPS_HAB
na_GVAgr	Gross Value Added growth	NUTS 2	PCH_PRE
pop_Total	Population	NUTS 2	NR
pop_Y[age]	Population by age	NUTS 2	PC_POP
training	Participation in education and training	NUTS 2	PC_YGE25-LE64



Methods: details



- Network Analysis

distance: weighted Euclidean distance

transformation: min-max transformation

edge threshold = {0.65, 0.8, 0.9}

algorithm: “Fruchterman–Reingold” algorithm

- Partition-around-medoids (PAM)

distance: weighted Euclidean distance

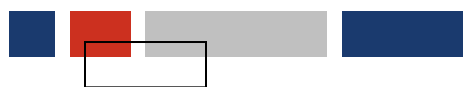
k -best = max(average silhouette width), when $k \neq 2$

- Over-representation analysis (ORA)

num2cat: $x \in Q_1$ and $x \in Q_3$

p-value correction: none

p-value cut-off = 0.10



Methods: details



■ Multinomial regression

remove NAs: column-wise ($NA > 0.00$) and row-wise ($NA > 0.00$)

remove predictors: i) r -between < 0.90 ; ii) r -within $< [\text{up-to } 30 \text{ vars}]$

transformation: min-max transformation

max terms: hard threshold (i.e. $nvars = nsubj - 1$)

confidence set = 100

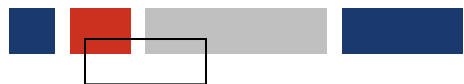
model level: only main effects

information criteria: AIC

models exploration: if $nmods < 200000$ exhaustive screening, else genetic algorithm

genetic algorithm: i) $popsiz = 100$, $mutrate = 10^{-3}$, $sexrate = 0.1$, $imm = 0.3$, $\delta M = 0.05$, $\delta B = 0.05$, $conseq = 5$; ii) $popsiz = 200$, $mutrate = 10^{-2}$, $sexrate = 0.2$, $imm = 0.6$, $\delta M = 0.005$, $\delta B = 0.005$, $conseq = 10$. Number replicates = 2.

model: multinomial log-linear model via single-layer feed-forward neural networks



Methods: details



■ Multivariate regression

remove NAs: column-wise ($NA > 0.00$) and row-wise ($NA > 0.00$)

remove predictors: i) r -between < 0.90 ; ii) r -within $< [\text{up-to } 30 \text{ vars}]$

transformation: min-max transformation

max terms: hard threshold (i.e. $nvars = nsubj - 3$)

confidence set = 100

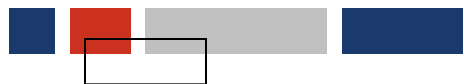
model level: only main effects

information criteria: AICc

models exploration: if $nmods < 200000$ exhaustive screening, else genetic algorithm

genetic algorithm: i) $popsiz = 100$, $mutrate = 10^{-3}$, $sexrate = 0.1$, $imm = 0.3$, $\delta M = 0.05$, $\delta B = 0.05$, $conseq = 5$; ii) $popsiz = 200$, $mutrate = 10^{-2}$, $sexrate = 0.2$, $imm = 0.6$, $\delta M = 0.005$, $\delta B = 0.005$, $conseq = 10$. Number replicates = 2.

model: multivariate linear regression



Methods: details



- Weighted correlation network analysis (WCNA)

remove NAs: row-wise (NAs > 0.50) and col-wise (NAs > 0.15)

distance = Topological Overlap Matrix($|Spearman\ r|^k$)

k -power transformation: $k = \min(k)$, when $k > 0.7 * k\text{-best}$

minimum module size = 2

cluster splitting level = 4, where level $\in \{1, 2, 3, 4\}$

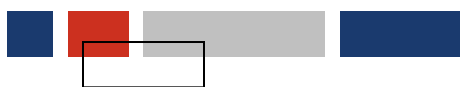
dynamic tree cut method = “hybrid”

merge cut-off height = 0.05

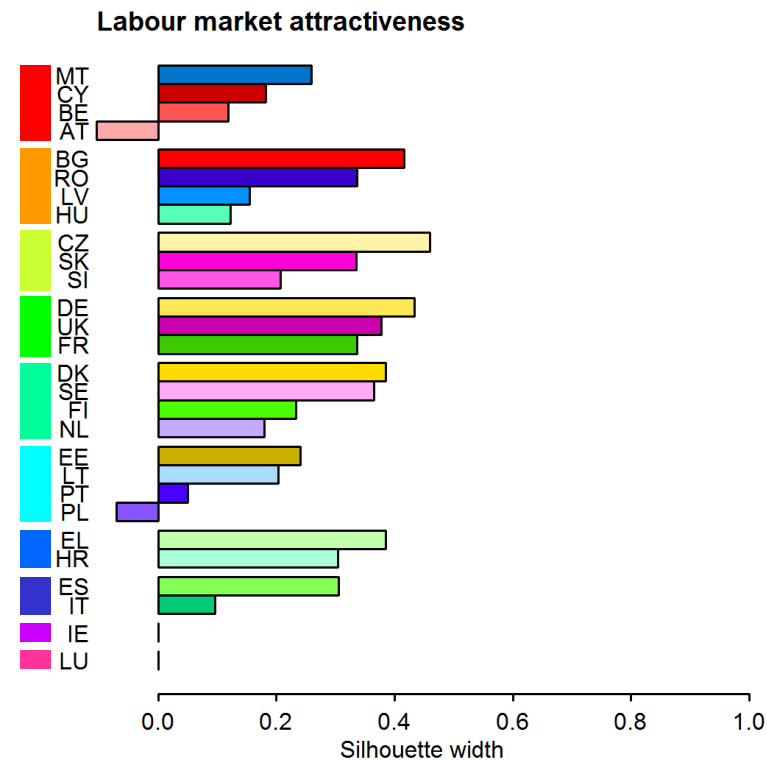
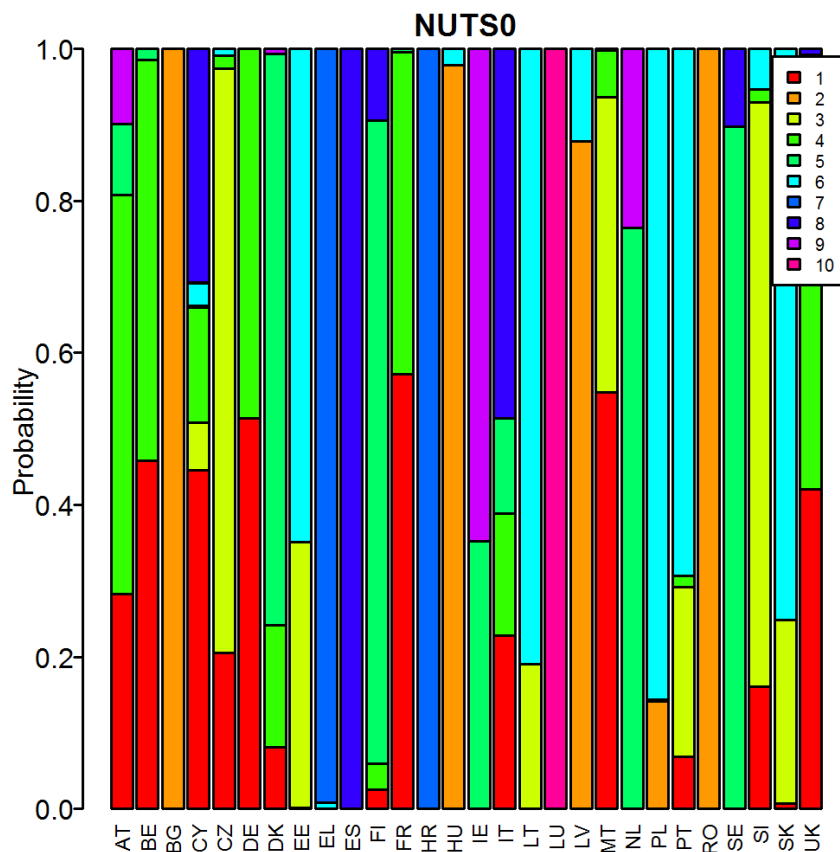
distance between modules = $1 - |Spearman\ r\ (\text{eigenvalues})|$

association modules and variables: i) Spearman r (eigenvalues) for numeric variables; ii) (R's McFadden)^{1/2} for categorical variables

p-value cut-off = 0.01



Results: LMkt Attractiveness



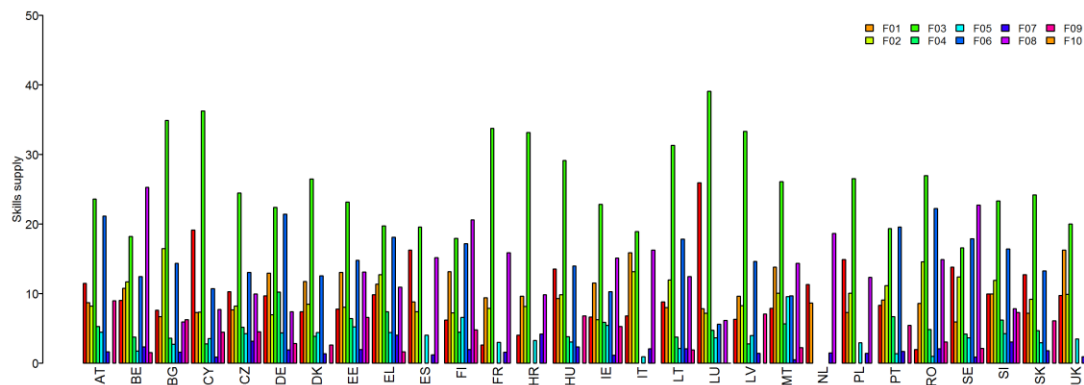
EU_groups ~ 1 + unemp_Y15-24 + GDP

$R_{\text{McFadden}} = 0.82$, $R_{\text{count}} = 0.82$



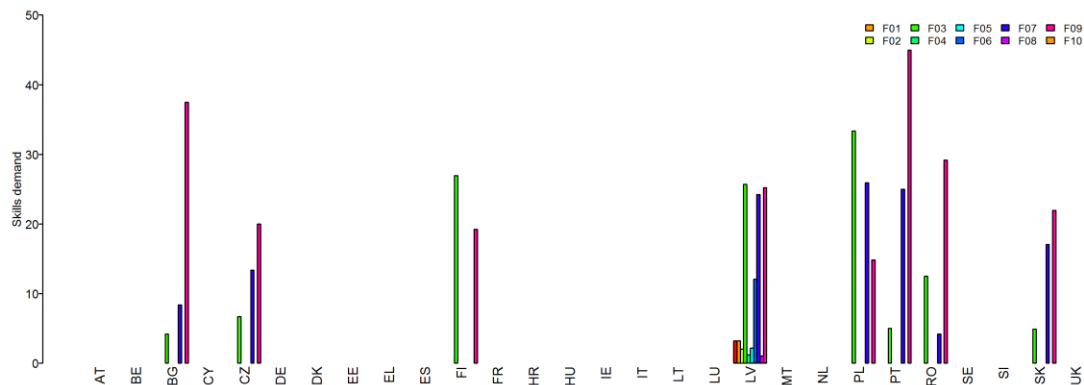


Results: Skills mismatch



Skills supply (“educ_uoe_grad02”):

Education fields (ISCED-F 13)



Skills demand (“jvs_a_nace2”):

Economic Activity (NACE Rev. 2)

Occupation Title (ISCO-08)

ad hoc mapping [... but ESCO v1]

distance metric (e.g. Maximum, Euclidean, Minkowski,...)

