

## PROBLEM STATEMENT

- TO BUILD A MODEL WHICH DEFINES THE PHYSICAL AND LOGICAL STRUCTURE OF TABLES
- THE MODEL IS USED TO DETECT TABLES AND TO ANALYZE AND DECOMPOSE DETECTED TABLES

## EXISTING SOLUTIONS

- CURRENT SOLUTIONS ARE BUILT AROUND PDFs
- THEY STRUGGLE WITH IMAGES AND ARE ALSO NOT EFFICIENT UNDERSTANDING THE UNDERLYING TABLE STRUCTURE

## CHALLENGE

BUILDING A MODEL WHICH DEALS WITH IMAGES AS EFFICIENTLY AS POSSIBLE  
WITH PDFs AND TO UNDERSTAND AND EXTRACT WIDE VARIOUS TABLE STRUCTURES

# COMPLEXITY

## NO BOUNDARIES

Your organisation's country of establishment (indicate your country of residence if answering as an individual person):

	replies	%
AT - Austria	79	(4.6%)
BE - Belgium	142	(8.2%)
DE - Germany	258	(14.9%)
DK - Denmark	28	(1.6%)
EL - Greece	80	(4.6%)
ES - Spain	123	(7.1%)
FI - Finland	22	(1.3%)
FR - France	132	(7.6%)
IE - Ireland	33	(1.9%)
IT - Italy	147	(8.5%)
LU - Luxembourg	1	(0.1%)
NL - Netherlands	89	(5.1%)
PT - Portugal	81	(4.7%)
SV - Sweden	68	(3.9%)
UK - United Kingdom	157	(9.1%)
CY - Cyprus	12	(0.7%)
CZ - Czech Republic	10	(0.6%)

**Table A Annual growth rates of debt securities issued by euro area**  
(percentages; end of period)

	2006	2007	2008
Total general government	2.5	2.8	8.1
<i>Long-term</i>	3.4	2.3	3.7
Fixed rate	3.4	2.0	3.5
Floating rate	3.4	5.4	5.1
<i>Short-term</i>	-8.8	9.6	61.9
Source: ECB.			

**Table B Structure of amounts outstanding of debt securities issued by euro area**  
(percentages of total debt securities issued by general government; end of period)

	2006	2007	2008
<i>Long-term</i>	92.9	92.5	88.8
Fixed rate	83.9	83.2	79.7
Floating rate	8.0	8.2	8.1
<i>Short-term</i>	7.1	7.5	11.2
Total general government in EUR billions	4,710.7	4,841.8	5,266.2
Source: ECB.			

# COMPLEXITY

## NESTED ROWS AND COLUMNS

**Table 6**  
**Correlations between the difficulties encountered and the perceived impacts on pupils, teachers and the school as a whole (Pearson's correlation coefficient\*)**

	Difficulties encountered with respect to		
	Lack of interest/ acceptance from colleagues	Lack of interest of pupils	Lack of interest of parents
<b>Impacts on participating pupils</b>			
Knowledge and awareness of different cultures	-0,1651	-0,2742	-0,1618
Foreign language competence	-0,0857	-0,1804	-0,1337
Social skills and abilities	-0,1237	-0,2328	-0,1473
Acquaintance of special knowledge	-0,1355	-0,2008	-0,1234
Self competence	-0,1291	-0,2466	-0,1636
<b>Impacts on participating teachers</b>			
Knowledge/appreciation of school system and education in the partner countries	-0,1505	-0,1636	-0,1349
Foreign language competence	-0,0545	-0,0997	-0,0519
Social skills and personal commitment	-0,2558	-0,2235	-0,1302
Professional knowledge and abilities	-0,2145	-0,2319	-0,1003
<b>Impacts on the school as a whole</b>			
European/international dimension of the school	-0,2438	-0,1945	-0,1030
School climate	-0,2976	-0,1810	-0,1012
Innovation in teaching and school management	-0,2586	-0,2557	-0,0928
Training of teachers	-0,1839	-0,1703	-0,0518
Involvement of external actors in the every day school-life	-0,2346	-0,2343	-0,2237
International mobility of pupils	**	**	-0,0583

\* Significance p = 0,000  
 \*\* No significant correlation

## NOISY IMAGES

Payslip for the month of May 2017

Employee No.:	P. D. S. H. R. G. L.	Location:	Department:
Name:			Designation:
Date of Joining:			Total Days:
Bank Name:			LOP:
Bank Acc. No.:			LOP Reversal:
PF No.:			Work Days:
ESI No.:			Effective Work Days:
PAN No.:			
UAN Number:			
<b>Earnings</b>			
Basic	14163	PF	
HRA	7081	INCOME TAX	
Special Allow	3982	PT	
Shift Allowance	2500		
Incentive	14131		
Statutory Bonus	1400		
<b>Total Earnings</b>	<b>43257</b>	<b>Total Deductions</b>	
Net Pay : Rs. 40280		(Rupees Forty Thousand Two Hundred Eighty)	

*D. S. H. R. G. L.*

# METHODOLOGY

PRE-  
PROCESSING

- SKEW-CORRECTION
- LINE EXTRACTION
- TEXT INVERSION

TABLE  
DETECTION

- BOUNDED FIGURES & TABLE DETECTION
- TEXT BLOCK FORMATION
- HORIZONTAL & VERTICAL ALIGNMENT VERIFICATION

TABLE  
EXTRACTION

- ROWS & COLUMNS IDENTIFICATION
- HEADERS' IDENTIFICATION
- MAPPING CELL'S ROW & COLUMN HEADER(S)

## INPUT

S12.5: Country Net Yearly Salary Average in EURO			
Country	Net Yearly Salary average	Country	Net Yearly Salary average
Austria	31.552	Italy	21.821
Belgium	25.408	Lithuania	9.150
Bulgaria	3.461	Luxembourg	6.340
Croatia	12.477	Malta	45.847
Cyprus	33.865	Netherlands	19.588
Czech Republic	11.514	Norway	36.791
Denmark	35.674	Poland	39.407
Estonia	7.965	Portugal	8.039
Finland	28.288	Romania	16.111
France	30.030	Slovakia	5.766
Germany	28.351	Slovenia	6.055
Greece	20.105	Spain	14.279
Hungary	9.423	Sweden	22.930
Iceland	33.936	Switzerland	28.075
Ireland	34.930	Turkey	64.123
Israel	26.703	United Kingdom	14.565
			37.001

Contract number - REM 01  
Final Report

Table 48 - Country Net Yearly Salary Average (2006; N=6.934, all currencies in EURO)

Note: A different analysis has been carried out for the net yearly salary costs obtained in the survey in order to detect unusual observations, the final sample had 7.018 correct replies, including 64 answers from Marie Curie fellowships. Marie Curie answers has not been considered for the calculation of the country net yearly salary average, as a result the sample had N=6.934.

## EXAMPLE

### [MAIN COLUMN HEADER/ROW HEADER/CO

	A	B	C
1	2	Country Austria	Net Yearly S
2	3	Country Belgium	Net Yearly S
3	4	Country Bulgaria	Net Yearly S
4	5	Country Croatia	Net Yearly S
5	6	Country Cyprus	Net Yearly S
6	7	Country Czech Republic	Net Yearly S
7	8	Country Denmark	Net Yearly S
8	9	Country Estonia	Net Yearly S
9	10	Country Finland	Net Yearly S
10	11	Country France	Net Yearly S
11	12	Country Germany	Net Yearly S
12	13	Country Greece	Net Yearly S
13	14	Country Hungary	Net Yearly S
14	15	Country Iceland	Net Yearly S
15	16	Country Ireland	Net Yearly S
16	17	Country Israel	Net Yearly S
17	18	Country Italy	Net Yearly S
18	19	Country Latvia	Net Yearly S
19	20	Country Lithuania	Net Yearly S
20	21	Country Luxembourg	Net Yearly S
21	22	Country Malta	Net Yearly S
22	23	Country Netherlands	Net Yearly S
23	24	Country Norway	Net Yearly S
24	25	Country Poland	Net Yearly S
25	26	Country Portugal	Net Yearly S
26	27	Country Romania	Net Yearly S
27	28	Country Slovakia	Net Yearly S
28	29	Country Slovenia	Net Yearly S
29	30	Country Spain	Net Yearly S
30	31	Country Sweden	Net Yearly S
31	32	Country Switzerland	Net Yearly S
32	33	Country Turkey	Net Yearly S
33	34	Country United Kingdom	Net Yearly S
34	35		

# PRE-PROCESSING

## SKEW-CORRECTION

## LINE DETECTION

## LINE REMOVAL

## NOISE REDUCTION

Caption-number: #103-1 Report-type: Net-Salary		
o S12.5: Country Net Yearly Salary Average in EURO		
Country	Net Yearly Salary average	Net Yearly Salary
Austria	25.400	25.400
Bulgaria	3.931	3.931
Croatia	12.150	12.150
Cyprus	13.665	13.665
Czech Republic	19.618	19.618
Estonia	16.731	16.731
Finland	35.071	35.071
France	30.120	30.120
Greece	23.357	23.357
Hungary	11.769	11.769
Iceland	9.463	9.463
Ireland	23.036	23.036
Italy	14.595	14.595
Netherlands	25.432	25.432
Portugal	7.885	7.885
Romania	10.111	10.111
Slovakia	8.055	8.055
Slovenia	11.779	11.779
Spain	12.310	12.310
Sweden	28.075	28.075
Turkey	12.423	12.423
United Kingdom	26.703	26.703
Total	37.051	37.051

Caption-number: #103-11 Report-type: Net-Salary		
o S12.5: Country Net Yearly Salary Average in EURO		
Country	Net Yearly Salary	Net Yearly Salary
Austria	25.400	25.400
Bulgaria	3.931	3.931
Croatia	12.150	12.150
Cyprus	13.665	13.665
Czech Republic	19.618	19.618
Estonia	16.731	16.731
Finland	35.071	35.071
France	30.120	30.120
Greece	23.357	23.357
Hungary	11.769	11.769
Iceland	9.463	9.463
Ireland	23.036	23.036
Italy	14.595	14.595
Netherlands	25.432	25.432
Portugal	7.885	7.885
Romania	10.111	10.111
Slovakia	8.055	8.055
Slovenia	11.779	11.779
Spain	12.310	12.310
Sweden	28.075	28.075
Turkey	12.423	12.423
United Kingdom	26.703	26.703
Total	37.051	37.051

Caption-number: #103-1 Report-type: Net-Salary		
o S12.5: Country Net Yearly Salary Average in EURO		
Country	Net Yearly Salary	Net Yearly Salary
Austria	25.400	25.400
Bulgaria	3.931	3.931
Croatia	12.150	12.150
Cyprus	13.665	13.665
Czech Republic	19.618	19.618
Estonia	16.731	16.731
Finland	35.071	35.071
France	30.120	30.120
Greece	23.357	23.357
Hungary	11.769	11.769
Iceland	9.463	9.463
Ireland	23.036	23.036
Italy	14.595	14.595
Netherlands	25.432	25.432
Portugal	7.885	7.885
Romania	10.111	10.111
Slovakia	8.055	8.055
Slovenia	11.779	11.779
Spain	12.310	12.310
Sweden	28.075	28.075
Turkey	12.423	12.423
United Kingdom	26.703	26.703
Total	37.051	37.051

Note: A different analysis has been carried out for the net yearly salary costs data and in the survey in order to detect unusual observations. The final sample obtained in the survey is used for the detection of outliers. The final sample had 70.8 correct answers, including 61 answers from Non-Career. Note that the mean salary average is 37.051, while the median salary average is 35.071. As a result the sample had N=6.934.

Note: A different analysis has been carried out for the net yearly salary costs data and in the survey in order to detect unusual observations. The final sample obtained in the survey is used for the detection of outliers. The final sample had 70.8 correct answers, including 61 answers from Non-Career. Note that the mean salary average is 37.051, while the median salary average is 35.071. As a result the sample had N=6.934.

Note: A different analysis has been carried out for the net yearly salary costs data and in the survey in order to detect unusual observations. The final sample obtained in the survey is used for the detection of outliers. The final sample had 70.8 correct answers, including 61 answers from Non-Career. Note that the mean salary average is 37.051, while the median salary average is 35.071. As a result the sample had N=6.934.

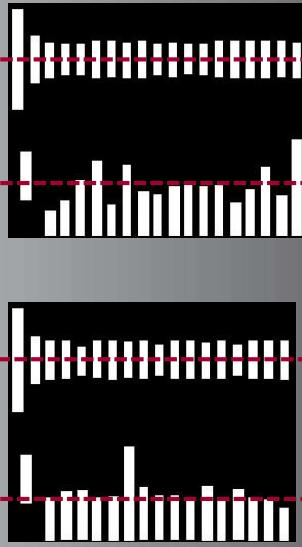
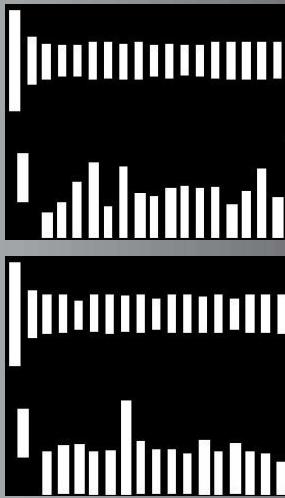
# TABLE DETECTION

BOUNDED FIGURES &  
TABLE DETECTION

BLOCK FORMATION

HORIZONTAL &  
VERTICAL ALIGNMENT  
VERIFICATION

Country	Net Yearly Salary average
Austria	31.552
Belgium	25.408
Croatia	3.361
Bulgaria	12.477
Cyprus	33.865
Czech Republic	11.514
Denmark	35.874
Estonia	7.865
Finland	26.288
France	30.030
Germany	28.351
Greece	20.105
Hungary	9.423
Iceland	33.936
Ireland	34.930
Israel	26.703



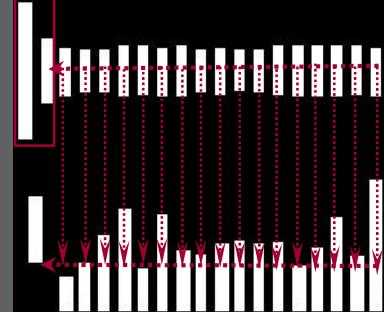
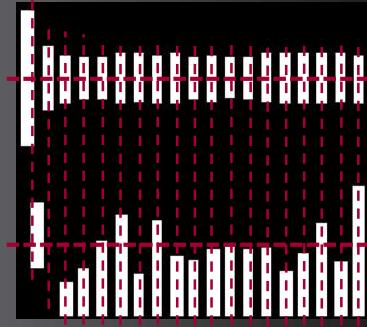
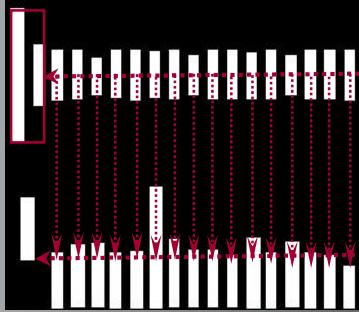
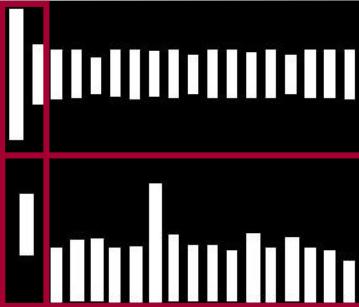
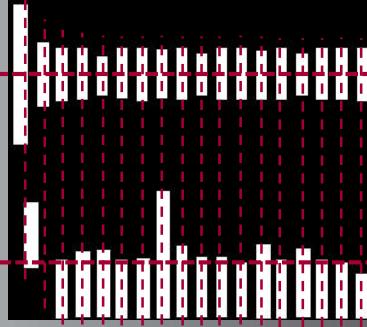
© 2023. Created and developed by DataFusion, Inc. All rights reserved.  
DataFusion is a trademark of DataFusion, Inc. All other trademarks and registered trademarks are the property of their respective owners. A different company from DataFusion, Inc. may own the rights to the trademarks and registered trademarks listed here. Some terms and conditions apply.  
This document is a work in progress and is subject to change without notice or obligation. The software, code, and documentation contained herein are provided "as is" and without warranty of any kind, either express or implied. The software, code, and documentation are provided for informational purposes only. The user agrees to use the software, code, and documentation at his/her own risk. The user agrees to indemnify and hold harmless DataFusion, Inc. from any claims, damages, losses, expenses, or costs arising out of the user's use of the software, code, and documentation.

## TABLE EXTRACTION

## ROWS & COLUMNS IDENTIFICATION

HEADER IDENTIFICATION

# MAPPING CELL'S ROW & COLUMN HEADER(S)



# CHALLENGES

## NESTED HEADER MAPPING

Scientific domain	Spain		Category with similar data from Eurostat	Total Yearly Salary Costs of Researchers (data from the study)	Category with similar p from
	Male	Female			
Social and Human Sciences	27.301 €	16.806 €	2.4.4	2.4	46.657 €
Economics	42.978 €	33.669 €	2.4.1	2.4	46.657 €
Chemistry	39.091 €	20.296 €	2.1.1	2.1	42.138 €
Physics	27.197 €	16.067 €	2.1.1	2.1	42.138 €
Life Sciences	36.523 €	22.139 €	2.2.1	2.2	37.111 €
Mathematics	45.389 €	41.107 €	2.1.2	2.1	42.138 €
Information Sciences	26.059 €	25.893 €	2.1.3	2.1	42.138 €
Engineering Sciences	34.316 €	25.435 €	2.1.4	2.1	42.138 €
Environment and Geosciences	29.210 €	14.847 €	2.1.4	2.1	42.138 €

# CHALLENGES

## REPETITIVE COLUMN/ROW HEADERS

Country	Number of responses	% of total responses	Country	Number of responses	% of total responses
DE - Germany	258	14.9%	CH - Switzerland	25	1.4%
UK - United Kingdom	157	9.1%	FI - Finland	22	1.3%
IT - Italy	147	8.5%	LT - Lithuania	22	1.3%
BE - Belgium	142	8.2%	IL - Israel	17	1.0%
FR - France	132	7.6%	PL - Poland	16	0.9%
ES - Spain	123	7.1%	CY - Cyprus	12	0.7%
NL - Netherlands	89	5.2%	CZ - Czech Republic	10	0.6%
PT - Portugal	81	4.7%	EE - Estonia	10	0.6%
EL - Greece	80	4.6%	HU - Hungary	10	0.6%
AT - Austria	79	4.6%	BG - Bulgaria	8	0.5%
SV - Sweden	68	3.9%	SK - Slovak Republic	7	0.4%
TR - Turkey	44	2.5%	SL - Slovenia	6	0.3%
IE - Ireland	33	1.9%	MT - Malta	5	0.3%
NO - Norway	31	1.8%	LV - Latvia	4	0.2%
Other country	30	1.7%	LU - Luxembourg	1	0.1%
RO - Romania	29	1.7%	IS - Iceland	1	0.1%
DK - Denmark	28	1.6%	<i>TOTAL</i>	1727	100%

# CHALLENGES

## TABLE SPLITTING OVER MULTIPLE PAGES

Guidelines for future European Union policy to support research		
1727 responses received <sup>1</sup> . 30/07/2004 - 15/10/2004		
<b>Details</b>		
I am answering as:		
Individual person	repplies	%
governmental body	540	(21.70)
university/higher education	141	(8.23)
commercial organisation (including consultancy) more than 250 employees	456	(26.45)
commercial organisation (including consultancy) less than 250 employees	115	(6.73)
association (e.g. trade associations, trade unions, employers' associations, chamber of commerce, NGO) or other (please specify)	144	(8.38)
Other (please specify)	113	(6.55)
Oneself	218	(12.63)
Your role in the organisation		
none - I am answering as an individual	350	(20.23)
senior management	350	(20.23)
researcher	180	(10.45)
strategic/policy function	619	(35.83)
specialist/expert	154	(8.90)
other (please specify)	278	(10.23)
210	(12.13)	
Your organisation's country of establishment (indicate your country of residence if answering as an individual person):		
AT - Austria	79	(4.85)
BE - Belgium	148	(8.23)
DE - Germany	258	(14.93)
DK - Denmark	28	(1.63)
EL - Greece	80	(4.63)
ES - Spain	123	(7.13)
FI - Finland	22	(1.36)
FR - France	132	(7.65)
IE - Ireland	33	(1.93)
IT - Italy	147	(8.55)
LU - Luxembourg	1	(0.05)
NL - Netherlands	89	(5.13)
PT - Portugal	81	(4.73)
SV - Sweden	68	(4.06)
UK - United Kingdom	157	(9.13)
CY - Cyprus	12	(0.73)
CZ - Czech Republic	10	(0.63)

Interactive Policy Making Online consultations		
EE - Estonia		
HU - Hungary		
LV - Latvia		
LT - Lithuania		
MT - Malta		
PL - Poland		
SK - Slovak Republic		
SL - Slovenia		
BG - Bulgaria		
TR - Turkey		
RO - Romania		
IS - Iceland		
LI - Liechtenstein		
NO - Norway		
CH - Switzerland		
IL - Israel		
Other:		
Your organisation's geographical area of activities (indicate your area of activities if answering as an individual person):		
repplies %		
local 57 (3.70)		
regional 116 (6.70)		
national 370 (21.45)		
European 226 (16.79)		
international 852 (49.35)		
not applicable 40 (2.36)		
Your organisation's activity type (indicate your activity type if answering as an individual person):		
repplies %		
Higher Education 402 (23.13)		
Research 818 (47.38)		
Industry 199 (11.53)		
Public Administration 80 (4.63)		
Other (please specify) 210 (13.35)		

<sup>1</sup> These are the overall figures based on all replies sent during the consultation. There are a couple of cases of coordinated research topic. In one case a group of researchers from different countries have put forward identical comments in the same form. However these two cases have not significantly affected the overall outcome of the statistics.

# CHALLENGES

## NOISY SCANNED DATA

Employee No	Bank Name
Name	Bank Acc. No.
Designation	IFSC Code
Original DOJ	PF No.
PAN	UAN
Payable Days	Paid Days

