A profile-aware methodological framework for collaborative multidimensional modeling

Sandro Bimonte INRAE, Clermont-Ferrand, France

sandro.bimonte@inrae.fr





Citizen observatories

Existing Scenario & Problems







Citizens collect data







Researchers work on data

Researchers define indicators and maps

INRAe





Volunteer data is « bad quality » data



Decision-makers, stake-holders and citizens INRA© do not understand results 2021

Citizen observatories

Our vision







Citizens collect data



Researchers work on data quality







Researchers & citizens co-design indicators



Researchers implement indicators (data and maps)







Researchers & citizens analyze indicators







Plan

- Background
- Case study
- Requirements
- Preliminaries
 - ICSOLAP & ProtOLAP
- Our new methodology
 - E-Pivot Table
 - Quality
 - Collaborative design
 - Experiments
- Conclusion & Future work





Background







Citizen science (1/2)

- Citizen science has been defined as an "online, distributed problem-solving and production model"
- Relies on crowdsourcing, i.e., on a large pool of people gathering inputs such as ideas, funding, etc.
- Volunteers:
 - provide data and
 - collaboratively edit them,

despite not being formally trained experts in the application domain







Citizen science (2/2)

- Citizen science has been applied in several application domains
 - Urbanism
 - Traffic jam
 - Ecology
 - Natural risks
 - ...
- Examples are crowdsourcing systems are numerous
 - OpenStreetMap, OpenEI, Wikipedia, etc.







Data Warehouse & OLAP (1/4)

 Data warehouse is: "A collection of integrated, subject oriented, time-variant, and non-volatile data from various sources into a single and consistent data repository to support decision-making" (Inmon, 1996)









- OLAP analysis
 - Drilling, Slicing, pivot, etc.
 - Interactive analysis and aggregation of huge volume of data
- OLAP : "A visual platform built especially to support rapid and easy spatiotemporal analysis and exploration of data following a multidimensional approach comprised of aggregation levels available in tabular and diagram displays"





Data Warehouse & OLAP (3/4)



✓ What is the total sold by month and city at 50 Km far from Paris







DW design process steps:





- Group Decision Support Systems (GDSSs): "Interactive computer-based environments that support concerted and coordinated team effort toward completion of joint tasks"
- Group facilitation is defined as a process in which a person (*facilitator*), who is considered as trustworthy by all the group members, intervenes to help improve the way they identify and solve problems, and make decisions.





Case study

- Bimonte, S. et al. VGI users & data centered methods for the analysis of farmland biodiversity indicators open issues. AGILE 2018, AGILE 2018 Lund, June 12-15, 2018
- Bimonte, S. et al. Collect and analysis of agro-biodiversity data in a participative context: A business intelligence framework. Ecol. Informatics 61: 101231 (2021)
- VGI4Bio.fr







Case study (1/3)

Project ANR VGI4Bio (www.vgi4bio.fr)

- mobilize two volunteers databases (Visionature and Observatoire Agricole de la Biodiversité - OAB) to build OLAP applications to analyze farmland biodiversity indicators.
- Visionature and OAB have 7682 and 1500 volunteers that produce data, respectively.

- Possible users interested in analyzing these data:

- the volunteers themselves that are interested in analyzing data to improve their data production quality, their related daily practices, etc.;
- public and private organisms (DREAL, Chambre d'Agriculture, etc.).
- These users are involved in the design of cubes
 - Unskilled OLAP, ICT users







Case study (2/3)

Example of data collection for pollinators





Sandro Bimonte



Soil



Resin







Mashed potatoes







Case study (3/3)

Example of OLAP analysis

departement	commune	(AII)	species	SOMME(abondance)	
BAS-RHIN			ecies	186	
			COTON	7	
			PETALES	0	
	BLAESHEIM-CODE-IN SEE:67049	All sp	ecies	75	
			COTON	0	
			PETALES	0	
	BOLSENHEIM-CODE-INSEE:67054	All sp	ecies	8	
			COTON	4	
			PETALES	0	
	DUNTZENHEIM-CODE-INSEE:67107	All sp	ecies	10	
			COTON	3	
			PETALES	0	
HAUT-RHIN		All species		21	
			COTON	0	
			PETALES	0	
	LAPOUTROIE-CODE-INSEE:68175	All species		6	
			COTON	0	
			PETALES	0	
	ROUFFACH-CODE-INSEE:68287	All sp	ecies	0	
			COTON	0	
			PETALES	0	







Requirements

Sakka, A., et al. **Volunteer Data Warehouse: State of Art.** International Journal of Data Warehousing and Mining (to appear)







Real life OLAP applications are quite complex

- Complex hierarchies
- Facts
- Aggregations





Requirements

Example of complex hierarchies





Requirement: Support complex Data Warehouse models







- Citizen science introduces another type of actors in DW projects: volunteers
- They are numerous and different from employed actors of classical DW projects





Requirements

• From our experience

Criteria	Involved users				
	Volunteers	Employed			
Knowledge of DW fundamentals	None / Very low	Low / Medium			
Involvement in the overall BI project	Partial	Full			
Geographical distribution	Very high	Very low			
Understanding of the project goals	Low / Medium	Medium / High			
Possibility of reaching unified vision	Very low	High			
Proficiency in the subject matter	Medium	High / Very high			
Availability to elicitation sessions	Low	High / Very high			
Number	Very high	Very low			







Volunteers define DW models that can present

• Similarities, Differences, and Conflicts





- Handling *divergent requirements*
- Handling requirements rejection

=> *Requirement:* Provide an additional design step with more engaged volunteers (i.e. *committers*) to finalize missing elements and solve requirements conflicts







- For volunteers is very difficult to validate the quality (completeness, minimality, relevance, etc.) of the provided DW schemata
 - Example: is the schema complete?, are some dimensions missing (such as for example land use, or urban

network)?







Requirement: the methodology must:

- be accompanied by some semi-structured interviews that allow guiding volunteers to identify the multidimensional schema problems in terms of semantic quality (i.e. quality based design)
 - The quality associated with data (such as distributive, normalization, etc.) is let to the DW experts since it depends on the structures of the DW
- be supported with some tools that allow correct (modify, delete, add) the DW schema elements that appears erroneous (i.e. *edit actions*)
 - Example: the temporal level 'week' could be removed since it is not explicative of any natural phenomena and is not useful for the analysis







- Due to the huge number of volunteers that might have many different needs in this project, providing an implementation for each proposed model is unrealistic because of its high human, temporal and financial costs
- Requirement: Reduce as much as possible time process from elicitation time to the prototypes (i.e. early rapid prototyping)
 - Example: in our case study we have 15 volunteers and each volunteer produce 3-5 DW schema in a very time expensive process







- Requirement: Usage of a very simple elicitation formalism because of volunteers have none or very low knowledge of DW fundamentals
 - Example: in our case study, some OAB volunteers are not comfortable also with uploading their data using a simple web page. They cannot understand complex design decision-making systems







							-			
	Handling	Handling	Quality	Correction	Prioritization	Committers	Early Rapid	Simple	Distributed	Design
	Divergent	Requirements'	Based	Actions		Involvement	Prototyping	elicitation	Time/Space	Approach
	Requirements	Rejection	Elicitation							
Bonifati	Manually	-	-	-	-	-	-	Interviews	-	Hybrid
et al	,									
2001										
WINTER			-	-	Ves	-	-	-	-	Beo-
ET										driven
41 2003										Griven
R.2005	Deview	n						International Action		
PAIMET	Review	Review	-	-	-	-	-	Interviews,	-	-
AL. 2003	sessions,	sessions,						Workshops,		
	Prototyping	Prototyping						Scenarios		
NABLI ET	-	-	-	-	-	-	-	2D sheets	-	Req-
AL. 2005										driven
GUO ET	-	-	-	-	-	-	-	Interviews	-	Hybrid
AL.										
(2006)										
GAM ET	Map	-	-	Yes	-	-	-	Map Formalism	-	Req-
AL (2007)	formalism									driven
GIORGINI	-	Basic	-	-	-	-	-	Interviews,	-	Reg-
ET AL.		operation						TROPOS		driven/
2005-		DW-tool								Hybrid
2008										
PRAKASH	-		-	-	-	-	-	GDI model	-	-
ET AL										
2008										
LUCIE ET	-		-	-	-	-	-	Interviews	Ver	-
JOCIKEI	-	-	-	-	-	-	-	Ourseline ws,	165	-
(2010)								Questionnaires,		
(2010)								reedbacks		-
ROMERO	-		-	-	-	-	-	rittering	-	Keq-
ET								functions		driven
AL.2010										
CRAVERO	-	-	-	-	-	-	-	-	-	Req-
ET										driven
AL(2013)										
KHOURI	Semantic	-	-	-	-	-	-	-	-	Req-
ET AL.	ontologies,									driven
2014	Pivot model									
KUMAR ET	Review	-	-	-	-	-	-	Interviews,	-	-
AL.(2014)	sessions							Workshops,		
								Prototyping,		
								Use cases, GDI,		
								DWARF		
DI-TRIA ET	-	-	-	-	-	-	-	-	-	Hybrid
AL(2015)										
ELAMIN ET	-	-	-	Yes	-	-	-	NI	-	Been
AL (2017)										driver
NASIBI ET	-	-	-	-	-	-	-	Guideliner	-	unven
AL 2017		-				-		autoennes		
AL. 2017	Famantic									0
RENET	Semantic		-	-		-	-	-	-	Keq-
AL.	ontologies									driven
(2018)	1	1	1	1	1	1	1			1



Preliminaries

ICSOLAP & ProtOLAP





ICSOLAP

•Boulil, K. et al. Conceptual model for spatial data cubes: A UML profile and its automatic implementation. Comput. Stand. Interfaces 38: 113-132 (2015)

•https://www.youtube.com/watch?v=2VQUrmU1yYk

•Available as Opensource for MagicDraw and Eclipse CASE tools






- To support advanced conceptual modeling we use in our framework the ICSOLAP UML Profile
- Our UML profile allows modeling advanced multidimensional issues:
 - complex dimensions (non-covering, non-onto and multiple hierarchies),
 - spatial dimensions,
 - multigranular facts,
 - complex aggregations, etc.
- We have successfully used ICSOLAP in several real applications: energy consuming, water quality, biodiversity







- We have implemented our profile in the commercial CASE tool MagicDraw
- OCL constraints are checked at design time







VGI4bio DW schema for the analysis of biodiversity in agricultural plots







ProtOLAP

- Bimonte, S. et al. **ProtOLAP: rapid OLAP prototyping with on-demand data supply**. DOLAP 2013: 61-66
- Bimonte, S. et al.Volunteered Multidimensional Design to the Test: The Farmland Biodiversity VGI4Bio Project's Experiment. DOLAP 2019
- <u>https://www.youtube.com/watch?v=WNHEwi4e3bg</u>
- Available as Opensource







ProtOLAP: methodology and tool for rapid DW prototyping







ProtOLAP tool (1/2)

- ProtOLAP takes as input the XMI file representing the ICSOLAP conceptual schema
- It automatically generates:
 - 1. the relational schema
 - 2. Mondrian XML metadata representing multidimensional elements, and the MDX-based calculated members for complex indicators
- This phase is completely automated without any intervention of designers → this accelerates software development and promotes standard processes



VARCHAR2(255)

NAME

alter table PRODUCTION add constraint CON 1

foreign key (A_DAYID) references DAY(DAYID).



<!-- The Mondrian definition of the EDEN indicators -->
<Measure name="sum" column="QUANTITY" aggregator="sum" visible="false" formatString="Standard" />
- <CalculatedMember name="AVG-SUM" dimension="Measures" visible="true" formatString="#,###.##">
<formula>Max{Descendants{[COOPERATIVES].CurrentMember, [COOPERATIVES.NodesHierarchy].[FARM]),[Measures].
 [sum])</formula>
<//CalculatedMember>
<//Cube>



ProtOLAP tool (2/2)

- A visual interface to feed both dimensions members and facts
- Example of dimensions feeding:
 - automatically create dummy members to saving time (e.g. temporal dimension)
 - manually by letting decision-makers insert the name of members to provide a better understanding (e.g. product dimension)







Our new design methodology

- Sakka, A. et al. A profile-aware methodological framework for collaborative multidimensional modeling. Data Knowl. Eng. 131-132: 101875 (2021)
- Bimonte, et al. Requirements-driven data warehouse design based on enhanced pivot tables. Requir. Eng. 26(1): 43-65 (2021)
- Sakka, A. A Volunteer Design Methodology of Data Warehouses. ER 2018: 286-300





DW Users Classification (1/2)

- We classify users based on their authoritativeness:
 - (non-authoritative) end-users: just express analysis requirements,
 - (authoritative) decision-makers: validate the end-users
 requirements and integrate them with their own requirements





DW Users Classification (2/2)



INRA











Quality-based volunteer collaborative design methodology









Quality-based volunteer collaborative design methodology





Quality-based volunteer collaborative design methodology







E-Pivot table requirement elicitation methodology

- Bimonte, S., L. Antonelli, S., Rizzi. **Requirement-Driven Data Warehouse Design Based on Enhanced Pivot Tables.** Journal Requirements Engineering (To appear)
- <u>https://www.youtube.com/watch?v=dQ5v8PNTiNY</u>











 Classical Pivot table do not support complex hierarchies

		· · · · · · · · · · · · · · · · · · ·
		Measures
species	type_traitement	MOYENNE(SOMME(MAX(abondance)POUR:annee)POUR:parcelle)
-All species	All type_traitements	1.727
PETALES	All type_traitements	0.098
RESINE	All type_traitements	0.108
COTON	All type_traitements	0.491
HERBES	All type_traitements	1.106
MORCEAUX-DE-FEUILLES	All type_traitements	1.273
FEUILLES-MACHEES	All type_traitements	2.069
TERREBOUE	-All type_traitements	6.943
	+MOLLUCIDE	5.876
	*FONGICIDE	6.592
	+AUTRE	6.653
	*HERBICIDE	6.789
	+TOUT TRAITEMENT	6.801
	+INSECTICIDE	6.902
	+PAS DE INSECTICIDE	6.962
	◆PAS DE AUTRE	6.991
	+PAS DE HERBICIDE	7.086
	+PAS DE TRAITEMENT	7.142
	+PAS DE MOLLUCIDE	7.179
	+PAS DE FONGICIDE	7.231



Slicor





E-Pivot table based methodology

- Pivot tables are queries over DW
 - Can be formally translated into DW models (Nabli et al., 2005)
 - Represent
 - Schema
 - Data
 - Aggregation

Decision-makers are familiar with pivot tables





E-Pivot table methodology: steps (1/4)

- 0. *Tutorial*. Presenting some existing OLAP applications to the decision makers, and explains to them the main concepts of DW and OLAP
- 1. Goal modeling (input: interviews; output: goal model)

Creating a model that defines the analysis goals and subgoals of each decision maker (Giorgini et al., 2008)





• Goal modeling example







E-Pivot table methodology: steps (2/4)

 2. E-Pivot table modeling (input: goal model, semistructured interviews; output: e-pivot tables model)

Refining the requirements specified in the goal models previously created

- First, decision makers express detailed requirements for each goal/subgoal by drawing e-pivot tables.
- Basically, e-pivot tables enhance classical pivot tables by establishing a graphical convention to visualize data in irregular hierarchies.
- Then, DW experts use a semi-structured interview to ask some specific questions to decision makers about their e-pivot tables in order to understand their needs better, and to incite decision makers to reason about possible errors and changes to be made.





E-Pivot table based methodology

- Examples of complex structures
 - strict hierarchy (a)
 - non-strict hierarchy (b)
 - non-onto hierarchy (c)
 - non-covering hierarchy (d)
 - many-to-many fact-dimension relationship (e)
 - multigranular fact (f)

Loca	Location				Tim	Time					Avg(abu	nda	nce)	
depa	artmen	nt	city			yea	year			month				
Haut	t-Rhin		Lapoutroie				2018			S	ep-18			13
Haut	t-Rhin		Rouffach				2018			C	Oct-18			12
Bas-	Bas-Rhin Strasburg			g		2018			Oct-18				14	
							(a	ı)						
	Locat	ion					Time Avg(a			oundance	2)			
	city			far	arm			year						
	Aubie	ubiere			Happy Farm			2017					13	
	Clerm	lermont-Fd												
	Aubiere			Bio	Bio Farm			2017			12			
	Aubiere			Bio	Bio Farm			2018					14	
							(1))						
	Boundary					Location				Avg(abundance)				
	type			flo	flower			plot						
	Wall						W	West Plot					12	
	Hedge			Rose			Ea	East Plot					15	
	Hedge			Jasmin			N	North Plot					13	
							(0	:)						
Crop						Time			ne		Avg(abu	Inda	ance)	
group		รเ	ubgro	up		crop	rop r			nth				
Field c	rop					Potate	otato				Sep-18			12
Field c	rop	G	rain			Corn					Oct-18			15
							(0	1)						
	Treatment					Tim			A	vg(abu	ndance)]		
	type produ Fungicide Prosa Fungicide Joao Fungicide Joao				product			year						
				aro		2018		3		12	1			
				0										
								7	13					
							(6)						
Species							1	Time Ave		g(abun	dance)			
oto	g	group spec			cies	s year								
ile	В	Butterfly					201)18			12		
	В	Butterfly Pie			ris		20	018		14				

(f)



E-Pivot table methodology: steps (4/4)

 4. Implementation (input: DW schema based on the ICSOLAP profile; output: SQL statements, Mondrian XML for DW prototype, and pivot tables issued from the DW prototype

For each ICSOLAP model, the DW experts implement a prototype of the DW using ProtOLAP and show it to the decision makers





Quality-based volunteer collaborative design methodology









Semi-automatic algorithm

















Quality: Filter

- Trustworthiness metric: assessing the reputation and reliability of end-users
 - It represents the degree of confidence that can be associated to each level and each indicator suggested by end-users
- Based on three factors: *Expertise*, *Reputation*, *Attractiveness*







Quality: Filter

- Expertise of end-user i (Expi)
 - Each end-user indicates her degree of expertise concerning the application domain of the cube schema using the ordinal scale {'low'; 'medium'; 'high'; 'expert'}







- Reputation of volunteer i. Reputation belongs to the community, not to the person whose trust was evaluated and it depends on many factors, including the previous behavior of that person
- Repi = f(1 rdi/di) where:
 - *rdi* represents the number of rejected data entries made by *i*,
 - *di* represents the number of all data entries entered by *i* (therefore *rdi* ≤ *di*);
 - f is a function with domain and codomain [0,..., 1]
 - In our case study $Repi = (1 rdi/di)^2$







- Attractiveness represents the "popularity" of a multidimensional element within the proposed cube schemata, and provides an important measure of its quality. Indeed, according to the many eyes principle, the more a piece of data is defined/edited, the higher its quality
- Attr_i(e) = {1, if user i has defined the element e
 0, otherwise}







Quality: Filter

• *Trustworthiness* of a multidimensional element *e* as: $Trust(e) = \sum i(w^{exp} \times Expi + w^{rep} \times Repi) \times Attri(e)/$ $\sum i(w^{exp} \times Expi + w^{rep} \times Repi)$ where: w^{exp} and w^{rep} are weights such that w^{exp} and $w^{rep} = 1$















- Completeness refers to sets of multidimensional elements, and indicates whether all necessary concepts have been modeled in the cube schema
- Precision refers to single multidimensional elements, and indicates whether an element is represented with sufficient precision and detail
 - Example: Time dimension is not precise enough since it starts at year level
- Relevance refers to single multidimensional element and indicates whether they are useful for analysis
- Minimality refers to sets of multidimensional elements and indicates whether they present redundancies







Endorse

- Consistency refers to single multidimensional element and indicates to what extent the rules characterizing the application domain have been adhered to
 - Example: for hierarchy group of species are grouped by specie
- Certainty refers to single multidimensional elements and indicates that there are no ambiguities in the names chosen for them
 - Example: Location used as name level for representing City is ambiguous







Endorse

- Confidentiality indicates whether a multidimensional element can be part of the cube schema or it should be removed for legal, anonymization, or confidentiality issues
- Usability indicates if, overall, a cube schema facilitates analysis







Quality attributes for cube schemata and their applicability to cube schema elements.

	Indicator	Set of indicators	Level	Set of levels	Hierarchy	Set of hierarchies	Dimension	Set of dimensions	Cube
Completeness		yes		yes		yes		yes	
Precision	yes						yes		
Relevance	yes		yes						
Minimality		yes		yes		yes		yes	
Consistency	yes				yes				
Certainty	yes		yes		yes		yes		
Confidentiality	yes		yes						
Usability									yes




















Actions on cube schem	na elements that can	be triggered in prese	ence of quality problen	ns during Endorse.
	Indicator	Level	Hierarchy	Dimension
Completeness	add	add	add	add
Precision	modify	add		
Relevance	delete	delete		
Minimality	delete	delete	delete	delete
Consistency	modify		modify	
Certainty	rename	rename	rename	rename
Confidentiality	delete	delete		
Usability	delete	delete	delete	delete





Methodology overview









 To take into account the varying users' profiles and skills, our framework recommends a different implementation of the Filter:

1. "*Are the end-users volunteers*?". This question is directed at the facilitator

- No: the reputation of all end-users is set to 1 when computing trustworthiness
- Yes: the reputation of end-users is computed

2. "Do you consider yourself an expert of the application domain?". This question is directed at the facilitator

- No: Filter step is automated based on the trustworthiness metric
- Yes: facilitator can manually select the well-defined elements of the cube schemata proposed by end-users













Collaboration: Endorse step

 To take into account the varying users' profiles and skills, our framework recommends a different implementation of the Endorse step depending on the answers give to the question:

"Do you think your skills in multidimensional modeling are good?". This question is directed at decision-makers.

- No for at least one of them: decision-makers must be supported by quality attributes and the *Endorse* step is based on majority votes
- Yes for all: considering that the evaluation of quality attributes can be long and tedious, the *Endorse* step can be carried out more informally via a free discussion





Collaboration: Implementation

 Vote and free discussion have been implemented in a webbased GDSS: GROUDA

Number of users that have submitted: 0 QualityBasedVote					
Element type	Element belongs to	Element to vote			
all dimensions	abundance	['location(60%)', 'species(30%)', 'Time(30%)', 'Farmer(30%)', 'Meteo(30%)', 'Crop(30%)']			
I don't want to participate.					
['location(60%)', 'species(30%)', 'Time(30%)', 'Farmer(30%)', 'Meteo(30%)', 'Crop(30%)']					
Completeness	Select your rank here		\$		
Minimality	Select your rank here		\$		
	Submit				





Experiments about quality metrics

- We have conducted 2 main experiments about quality metrics of Endorse step
- All other steps have been also validated with real experiments









Conclusion & future work







- Involve volunteers in the definition of analysis is crucial for make the citizen observatory survive
- Use OLAP to analyze volunteers data by volunteers
 - Yes but, how design a DW with such kind of non-skilled DW users?
- We propose a new methodology based on:
 - E-pivot table
 - Group decision support systems
 - Rapid prototyping for validation
 - Quality metrics
- Each step is implemented with an "available tool" (ask me ;-))
- Validated on real and academic case studies







- We are finalizing a web based tool for E-Pivot Table
- Automatize the data source validation





Thanks to my collegues and students that worked with me

Questions ?

sandro.bimonte@inrae.fr



