

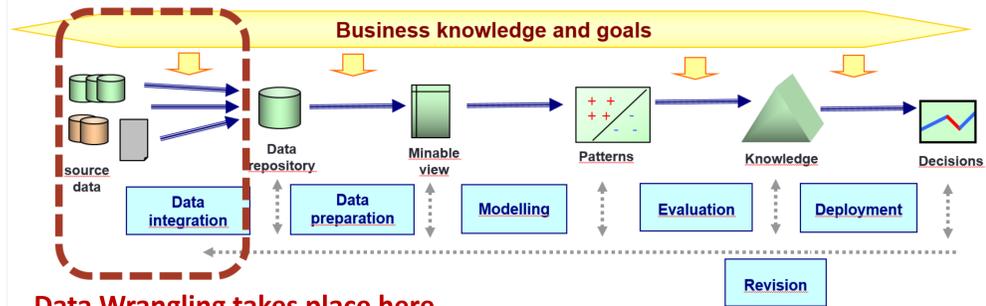
General-Purpose Inductive Programming for Data Wrangling Automation

Lidia Contreras-Ochando, Fernando Martínez-Plumed, Cèsar Ferri, José Hernández-Orallo and María José Ramírez-Quintana
Universitat Politècnica de València (UPV), Spain
{liconoc, fmartinez, cferri, jorallo, mramirez}@dsic.upv.es



UNIVERSITAT POLITÈCNICA DE VALÈNCIA

DATA WRANGLING IN DATA SCIENCE



Data Wrangling takes place here

- 50%-80% of effort.
- It depends on previous knowledge.

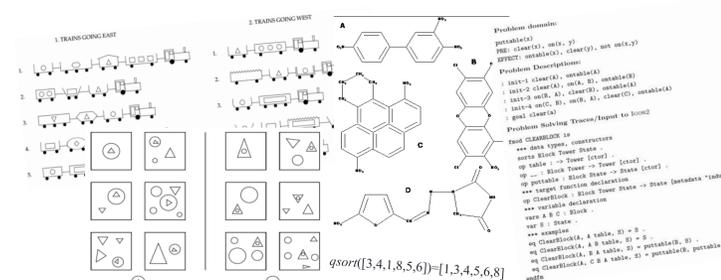
(Semi-)Automating Data Wrangling would have a very significant impact

INDUCTIVE PROGRAMMING (IP)

- Machines program themselves (programming by example, program synthesis, ILP, ...).

Input:

- Data D (usually small sets of data)
- Background Knowledge BK (facts, functions, constraints, etc.)



Output:

- Hypothesis h (a program, possibly rich and complex)

```

qsort([], []) ←
qsort([X|T], S) ← part(X, T, L1, L2),
                  qsort(L1, S1),
                  qsort(L2, S2),
                  app(S1, [X|S2], S)

delOdds(L, R) ← L=[_], R=[]
delOdds(L, R) ← L=[_|L1], odd(HL), delOdds(TL, TR), R=[_|TR]
delOdds(L, R) ← L=[_|L1], _=odd(HL), delOdds(TL, TR), R=[_|TR]

B :- Vx. ((group_likes x) =
         if (L1, B_likes B_likes x) then T
         else if (L1, 1:0: likes(X,Y), friendof(X,Y),
                0.8: likes(X,Y), friendof(X,Z), likes(Z,Y),
                0.5: friendof(john,mary),
                0.5: friendof(mary,pedro),
                0.5: friendof(pedro,tom).
         ...

```

IP FOR DATA WRANGLING: GPDL + BK!

- Re-use any IP system based on a **general-purpose declarative language (GPDL)**:
 - GOLEM, Progol, (F)FOIL, ADALGO, FLIP, IGOR I/II, Aleph, MagicHaskell, Metagol, gErl, ...
 - Users may edit the solutions written in languages such as Haskell, Erlang or Prolog.
- User can choose or modify the libraries and **background knowledge (BK)**

EXAMPLE

- Attribute transformation performed with:
 - **MagicHaskell** (Katayama 2004-today): <http://nautilus.cs.miyazaki-u.ac.jp/~skata/MagicHaskell.html>
 - **Metagol** (Muggleton et al. 2014-today) [.https://github.com/metagol/metagol](https://github.com/metagol/metagol)
 - **gErl** (Martinez-Plumed et al. 2012-today). <https://github.com/nandomp/gErl>

IP FOR DATA MANIPULATION: THE DSL APPROACH

- Data manipulation automation using DSLs
 - Define IP systems over Domain Specific Languages (DSLs) fitting the domain.
 - **Expressive enough** to cover the problems in the domain.
 - **Restrictive enough** to enable efficient search.
- It has been a success! (Gulwani 2011-2016)
 - FlashFill, FlashExtractText, FlashRelate, FlashNormalize, BlinkFill, ...
- **Limitations:**
 - Systems (IP engines) **must be redesigned** for each domain...
 - **Lack of flexibility:** customisation, including background knowledge, ...

Flash Fill:

	A	B
1	Email	Column 2
2	Nancy.FreeHafer@fourthcoffee.com	nancy freehafer
3	Andrew.Cencic@northwindtraders.com	andrew cencic
4	Jan.Kotas@litwareinc.com	jan kotas
5	Mariya.Sergienko@gradicdesigninstitute.com	mariya sergienko
6	Steven.Thorpe@northwindtraders.com	steven thorpe
7	Michael.Nelpper@northwindtraders.com	michael nelpper
8	Robert.Zare@northwindtraders.com	robert zare
9	Laura.Giussani@adventure-works.com	laura giussani
10	Anne.Hl@northwindtraders.com	anne hl
11	Alexander.David@contoso.com	alexander david
12	Kim.Shane@northwindtraders.com	kim shane
13	Manish.Chopra@northwindtraders.com	manish chopra
14	Gerwald.Oberleitner@northwindtraders.com	gerwald oberleitner
15	Amr.Zaki@northwindtraders.com	amr zaki
16	Yvonne.McKay@northwindtraders.com	yvonne mckay
17	Amanda.Pinto@northwindtraders.com	amanda pinto

Flash Extract:

Label 1	Label 2	Label 3
757) 555-8534	Redmond	(757) 555-1834
441) 555-5788	Redmond	(441) 555-5788
441) 555-5779	Seattle	(441) 555-5779
441) 555-2776	Redmond	(441) 555-6774
441) 555-5779	Seattle	(441) 555-2462
441) 555-5779	Redmond	(888) 555-2770

	Gender	GenderOK	Birthdate	BirthdateOk	Postcode	PostcodeOK	Score	ScoreOk	Km	metresOK	Weight	WeightOK
#1	Male	M	3 1 1971	1971 1 3	46 025	46025	5.5, 4.6, 5.8	5.3	5	5000	f "CAMP DRY DBL NDL 3.6 OZ"	"3.6 OZ"
#2	Female	F	4 5 1993	1993 5 4	46225	46225	3.5	3.5	3	3000	f "DRY NDL 0.23 KG"	"0.23 Kg"
#3

STATUS

- **Analyse common transformations** in data wrangling (from tutorials/books/users/systems and DSL-based wrangling tools).
 - Feature transformations.
 - Row transformations.
 - Table transformations.
 - Integration from several tables.
 - Other kinds of formatting.
- **Define a library (BK)** that *can* solve many of these common problems.
 - Trade-off between efficiency and power (syntactic and semantic domain).
- BK selection by **interaction** (asking the user)
 - Do you think this is a date? An address? Use the appropriate BK for the task.
- With MagicHaskell:
 - Identify the function library (LibTH.hs) to make it easily editable:
 - Modifiable by the user and coding new functions.
- With Metagol:
 - Try to resolve some type problems
 - Two "libraries" to be selected for Metagol: Metarules and proper BK.
- With gErl:
 - Needs learning operators, data representation and examples navigation.
 - Identify BIFs or new functions to be added in the background knowledge.

Download the paper:



<https://goo.gl/YvC75L>