

Adastra Bulgaria and Data Quality

July Karchev DWH Technical Lead

Georgi Pamukov, Galina Guneva DWH Consultants

October 12, 2011 CAREER DAYS 2011, IT, TELECOMMUNICATIONS, BPOSofia



Adastra Group



- Established in 1994, headquarter in Toronto, Canada
- Offices in 7 countries
- Over 420 successful projects in 24 countries
- Over 800 consultants
- 32 international quality awards



Adastra Bulgaria

- Business Information Management Leader
- 100+ IT professional consultants
- International and local projects for world-leading companies



3



- Sofia office
 - 29th Panayot Volov street

SADASTF

- Varna office
 - 2 Dunav street

Selected Partners





SADASTRA

Data Quality – the problem

Everyone is affected by poor data quality

- Directly 2 or 3 identical mailings from the same sales organization in the same week
- In less direct ways 20 minute wait on hold for a customer service department
- Dangerously deliberate identity theft
- Impact on business and government agencies
 - Ineffective operations- data quality problems cost companies about 10% of their revenue



Bad Data – You'll Cry

Low quality data costs billions

Gartner has reported

- 25% of critical data within large businesses is somehow inaccurate or incomplete
- and that 50% of implementations fail due to a lack of attention to data quality issues
- data quality problems cost U.S. businesses over US\$600 billion a year



What is Data Quality

Data quality is something everyone desires

- Everybody has a different idea of what data quality means
 - Each organization has own rules on the quality of data
 - Each company department has own data quality expectations
- Definition of data quality
 - "fitness to serve its purpose" to what extent the data are appropriate for their purpose
 - In practice this means identifying assets of data quality objectives associated with any dataset and then measuring that datasets conformance to these objectives



The many faces of Data Quality

Accuracy - does data reflect real world state

- company identifier exists in the official register of companies
- street, street number, city, ZIP, country code found in an authorized/official etalon of addresses
- Completeness extent to which the expected attributes of data are provided
 - Empty address fields, blank characters, "_", "-"
- Consistency does pieces of info match together
 - Трабант 2105
- Timeliness/Relevance is the data up-to-date
 - Old pricing info
- Validity is the data in a reasonable range/format

2090



Accuracy – data quality gone wrong

	src_name	src_gender	<pre>src_birth_date</pre>	src_sin	src_card	src_address
1	Dr. John Smith	М	12/16/1978	00000000	88682239496	14618 110 Ave Surrey V3R2A9
2	Smith W. John	MALE	16.12.1978	095-242-434	266805807984498	Surrey 14618 110 Ave
3	John William Smith		781216	SIN095242434		25 Linden Str Toronto M4X 1V5
4	Dr. J.W. Smith	M	11/16/78	095242433	4334338874158398	
5	John Smith		16.11.1978	095252433	43176816175728 So	cial Insurance Number
6	Smith John		16.11.1978		43176816175728	not corresponding
7	John Smiht		16.11.1978	95252433	NULL to au	uthorized/official etalons
8	Jane Watson		1982	420347213	974997714973	
9	Watson Jane	FEMALE	5.1.1982	420-347-213		8500 Leslei street Toronto L3T 7M8
10	Jane Smith	F	1982-05-01	SIN420347213	720847891758473	
11	J. Smith			420-347-213		
12	Janette Smith	F	B. 2015-01-05	SIN	54590972345186400	

Inconsistent pattern of values
in the field

There is no such numbers on this streets

	src_street	src_city	src_province	src_zip
4	8500 Leslie	Toronto	Ontario	L3T7M8
2	128 Yonge St	Tronto	ON	M5C1XX
3	25 Linden Str	TORONNTO	SOFIA OKRAG	2227
4	8500 Leslie street	Marham	Onntario	L37
5	8500 Leslei street	Toronto	Ontsrio	L3T 7M8
B	8500 Leslei st	Totonto	on	L3T 7M8
7	128 St John St	St John	NB	
8	20 St John Ave	Oromocto	NB	
9	5867 Eagle Island	Vancouver	Britisch Columbia	V7W1V5



Accuracy – data quality gone wrong (cont.)

			Dupl this is s record	icated data – several different is for the same person		
	src_name	src_gender	src_birth_date	src_sin	src_card	src_address
1	Dr. John Smith	М	12/16/1978	0000000	88682239496	14618 110 Ave Surrey V3R2A9
2	Smith W. John	MALE	16.12.1978	095-2 <mark>42-434</mark>	266805807984498	Surrey 14618 110 Ave
3	John William Smith		781216	SIN058242434		25 Linden Str Toronto M4X 1V5
4	Dr. J.W. Smith	М	11/16/78	095242433	4334338874158390	
5	John Smith		16.11.1978	095252433	431768161757282	8500 Leslie L3T 7M8 Toronto
6	Smith John		16.11.1978		431768161757282	8500 Leslie street Marham
7	John Smiht		16.11.1978	95252433	NULL	
8	Jane Watson		1982	420347213	974997714973	600-8500 Bugatica str. Toronto L3T 7M8
9	Watson Jane	FEMALE	5.1.1982	420-347-213		8500 Leslei street Toronto L3T 7M8
10	Jane Smith	F	1982-05-01	SIN420347213	720847891758473	
11	J. Smith			420-347-213		
12	Janette Smith	F	B. 2015-01-05	SIN	54590972345186400	



Completeness – data quality gone wrong

	src_name	src_gender	src_birth_date	src_sin	src_card	src_address
1	Dr. John Smith	М	12/16/1978	00000000	88682239496	14618 110 Ave Surrey V3R2A9
2	Smith W. John	MALE	16.12.1978	095-242-434	266805807984498	Surrey 14618 110 Ave
3	John William Smith			SIN095242434		25 Linden Str Toronto M4X 1V5
4	Dr. J.W. Smith	M Missin	a date & month	095242433	4334338874158390	
5	John Smith	MISSI	ig date a month	095252433	431768161757282	8500 Leslie L3T 7M8 Toronto
6	Smith John		16.11.1978		431768161757282	8500 Leslie street Marham
7	John Smiht		16.11.1978	95252433	NULL	
8	Jane Watson		1982	420347213	974997714975	600-8500 Bugatica str. Toronto L3T 7M8
9	Watson Jane	FEMALE	5.1.1982	420-347-213		8500 Leslei street Toronto L3T 7M8
10	Jane Smith	F	1982-05-01	SIN420347213	720847891758473	
11	J. Smith			428-347-213		Missing/blank Values
12	Janette Smith	F	B. 2015-01-05	SIN	54590972345186400	

		src_street	src_city	src_province	src_zip
'-' values in 'src_province' field	22	618 REVTELL SYND NW	EDMONTON	AB	T6R2H9
	23	283 PARKE STEEET S	HAMILTON	ON	L8P3G5
	24	2046 SIRWOCCO DVIRE SW	CALGARY	-	T3H2M8
	25	15 MCDGILL STRETE S	SMITHS FALLS	-	K7A3M4
	26	229 SPRUCEE STRET S	TIMMINS	-	P4N2M8
	27	4508 BEEDIE	BURNABY	-	V5J5L2
	28	1460 ESTHER-BLONDIN	QUEBEC	QC	G1Y3N7
	29	295 FATHER TOBIN	BRAMPTON	ON	L6R0N2
	30	273 RIVER	SUNNY CORNER	NB	E9E1C9

Consistency – data quality gone wrong





Validity – data quality gone wrong

	src_name	src_gender	src_birth_date	src_sin	src_card	src_address
1	Dr. John Smith	М	12/16/1978	00000000	88682239496	14618 110 Ave Surrey V3R2A9
2	Smith W. John	MALE	16.12.1978	095-242-434	266805807984498	Surrey 14618 110 Ave
3	John William Smith		781216	SIN095242434		25 Linden Str Toronto M4X 1V5
4	Dr. J.W. Smith	М	11/16/78	095242433	4334338874158390	
5	John Smith		16.11.1978	095252433	431768161757282	8500 Leslie L3T 7M8 Toronto
6	Smith John		16.11.1978		431768161757282	8500 Leslie street Marham
7	John Smiht		16.11.1978	95252433	NULL	
8	Jane Watson		1982	420347213	974997714973	600-8500 Bugatica str. Toronto L3T 7M8
9	Watson Jane	FEMALE	5.1.1982	420-347-213	K	8500 Leslei street Toronto L3T 7M8
10	Jane Smith	F	1982-05-01	SIN420347213	720847891758473	
11	J. Smith			420-347-213		
12	Janette Smith	F	B. 2015-01-05	SIN	54590972345186400	

Impossible date of birth – value out of range Not valid card numbers (too short)



Data Quality Market

Growing need for data quality analysts

- Employment of data quality analysts was projected to increase 20% from 2008 to 2018, according to the U.S. Bureau of Labor Statistics (BLS).
- According to a Gartner report, the BI and analytics market is expected to reach \$10.8 billion in 2011.



Data Quality – anchors

- Data Quality Experts and Consultants
 - Business Analyst
 - Data Analyst
 - Data Steward
- Methodologies
- Software Tools



Case Study – DQA Project

The Company

Company with 2 million customers, 3 million accounts. Enterprise information system handles customer registration and servicing, billing, invoicing and payments.

The Problem

Company reported difficulties with target marketing campaigns and inability to send invoices due to incorrect or missing address information.

Corporate reports have showed differences in sold products versus billed amounts.

The Challenge

Identify data with quality issues and implement solution for DQ improvement. Evaluate the amount of impacted address information and resolve as much as possible of the affected information.



Project environment: DQ Analyzer [PROFILE & ANALYZE]





Data Quality Assessment

Various data analyses to reveal basic and hidden DQ issues

- Unique data profiling tool to reveal basic and hidden data quality issues
- Easy to work with (tutorials and illustrative samples)
- Completely FREE for commercial use
- www.ataccama.com



Unmatched performance

Analyze millions of DB or CSV records in minutes using an easy-to-use wizard

	_	١
-	_	
_		٦
		J

Regular Expressions

Validate format/structure of the data, extract partial information from unstructured text



Rule-based Engine

Context-based validation using customizable rules



Project environment: Ataccama DQC [CLEANSE & MATCH]





Data Cleansing & Enrichment

Improve quality of individual records, enrich the data using external data sources.



Match & Merge

Correctly match related records, create representative "golden records".



DQ Firewall

Prevent poor quality data from entering the systems by leveraging DQC as the validation procedure.



DQ Reporting & Monitoring

Set up and run DQ reports periodically to monitor quality of your data.

- Essential tool for complex DQM designed to evaluate, monitor and manage quality of data
 - Assessment [Data Profiling]
 - Prevention [DQ Firewall]
 - Measurement [DQ Reporting]
 - Control [DQ Monitoring]
- Bundled with
 - Vertical-specific and country-specific sets of business rules
 - Localized dictionaries and knowledge bases
- Flexible and platform-independent
- Scalable and high-performance oriented (incremental batch and online mode)



Activities in a Data Quality Project

1. Profiling

basic analysis - metadata discovery and definition

2. Parse

extract individual elements and store in correct fields

3. Cleanse and Standardize

remove non-relevant information and "noise" from the content of the data reach uniform structure and enrich fro etalons

4. Match and Merge

identify and consolidate records that refer to the same business object(customer for example)

5. Enrich

adding useful, but optional, information to existing data or complete data



Project activities – 1. Profiling (Basic)

Basic Frequency Mask Quantiles Groups **Basic Analyses** What the data Expression: TELNO Data type: STRING Domain: Rows: 619,251 analysis revealed? Counts **Completeness issues: Accuracy issues: Missing values Potentially duplicated** Unique Null records 46.29% 46.49% Duplicate Non-unique 2.97% 4.26% **Accuracy issues: Non-standard values** % Туре Count Null 287,868 46.49% 331,383 53 51% Non-null Duplicate 26,369 4.26% **Accuracy issues:** Distinct 49.26% 305,014 Non-unique 18,392 2.97% Names in telephone Unique 286,622 46.29% **Accuracy issues:** number column **Inconsistent pattern** of values Statistics in the field Туре Value Frequency Minimum value 3 Median value 8243663 янчева 99 Maximum value Type Value Minimum length 1 Median length 9 Average length 9.19 Maximum length 20



Project activities - 1. Profiling (Mask)

What the data analysis revealed?

Accuracy issues: Inconsistent pattern of values in the telephone number field L – represents letter D – represents digit Basic Frequency Mask Quantiles Groups

Mask Analysis

Mask: characters: [:letter:] -> L,[:digit:] -> D

Value	Count	%	
LLLLL LL.DD DDDDDDD	1	0.00%	
LLLLL LLL.	1	0.00%	
	1	0.00%	
LLLLL LLLLL-DDDDDDD	1	0.00%	
LLLLLL-DDDDDD	1	0.00%	
LLLLLL-DDDDDDDD	1	0.00%	
LLLLLL-DDDDDDDDD	1	0.00%	
LLLLLL, DDDDDDD	1	0.00%	
LLLLLL.DDDDDDD	1	0.00%	
LLLLLL DDDDDD	1	0.00%	
LLLLLL - DDDDDDDDDD	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
111111-	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	
	1	0.00%	

Project activities - 2. Parsing

Goal:

- When different types of data are in a single field, extraction of individual elements and storing in correct fields are needed;
- This will also allow performing of cleansing, standardizing and enrichment of the data.

_						-	1				
			_				PARS	SING			
1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.ХРИ	СТО БЕЛЧЕВ/	ГАВРИЛ ГЕНОВ/	Crp:21		Splitting of th	e individual			
2 1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.ИВА	Н ДЕНКОГЛУ	/ Crp:44	<null></null>		olomonto a	ad storing			
3 1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.ХРИ	СТО БЕЛЧЕВ,	/ГАВРИЛ ГЕНОВ/	Crp:18		elements al	iu storing			
1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.АНГ	ЕЛ КЪНЧЕВ М	Vº2	<null></null>		in correc	t fields			
5 1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.ГРА	Ф ИГНАТИЕВ	Crp:12	<null></null>						
5 1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.ГРА	Ф ИГНАТИЕВ	Сгр:14	<null></null>						
/ 1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.ГРА	Ф ИГНАТИЕВ	Бл:16	<null></null>						
3 1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.ГРА	Ф ИГНАТИЕВ	Сгр:16 Ет:1	<null></null>						
9 1000 СРЕДЕЦ	Ц КВ.ЦЕНТЪР	УЛ.ГРА	Ф ИГНАТИЕВ	Сгр:16 Ет:2	<null></null>						
10 1000 CPE / FU	I KB.LIEHTЪP	УЛ.ГРА	Ф ИГНАТИЕВ	Nº 18	<null></null>						
10 1000 Cr Edita	,										
		TYPE	DISTRICT	STREET TYPE	STREET	STREET NO			ENTRANCE	FLOOR	
POSTAL_CODE 1	MUNICIPLITY	TYPE KB.	DISTRICT LEHTЪP	STREET_TYPE УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ	STREET_NO	BLOCK_NO <null></null>	BUILDING_NO	ENTRANCE <null></null>	FLOOR <null></null>	
POSTAL_CODE 1 1000 0 0	MUNICIPLITY СРЕДЕЦ СРЕДЕЦ	TYPE KB. KB.	DISTRICT LEHTЪP LEHTЪP	STREET_TYPE УЛ УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ ИВАН ДЕНКОГЛУ	STREET_NO <null> <null></null></null>	BLOCK_NO <null> <null></null></null>	BUILDING_NO 21 44	ENTRANCE <null> <null></null></null>	FLOOR <null> <null></null></null>	
POSTAL_CODE 1 1000 0 1000 0	MUNICIPLITY СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ	TYPE KB. KB. KB.	DISTRICT LEHTЪP LEHTЪP LEHTЪP	STREET_TYPE УЛ УЛ УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ ИВАН ДЕНКОГЛУ ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ	STREET_NO <null> <null> <null></null></null></null>	BLOCK_NO <null> <null> <null></null></null></null>	BUILDING_NO 21 44 18	ENTRANCE <null> <null> <null></null></null></null>	FLOOR <null> <null> <null></null></null></null>	
POSTAL_CODE 1 1000 0 1000 0 1000 0	MUNICIPLITY СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ	TYPE KB. KB. KB. KB.	DISTRICT ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР	STREET_TYPE УЛ УЛ УЛ УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ ИВАН ДЕНКОГЛУ ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ АНГЕЛ КЪНЧЕВ	STREET_NO <null> <null> <null> 2</null></null></null>	BLOCK_NO <null> <null> <null> <null></null></null></null></null>	BUILDING_NO 21 44 18 <null></null>	ENTRANCE <null> <null> <null> <null></null></null></null></null>	FLOOR <null> <null> <null></null></null></null>	
POSTAL_CODE 1 1000 0 1000 0 1000 0 1000 0	МUNICIPLITY СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ	TYPE KB. KB. KB. KB. KB.	DISTRICT ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР	STREET_TYPE УЛ УЛ УЛ УЛ УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ ИВАН ДЕНКОГЛУ ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ АНГЕЛ КЪНЧЕВ ГРАФ ИГНАТИЕВ	STREET_NO <null> <null> <null> 2 <null></null></null></null></null>	BLOCK_NO <null> <null> <null> <null> <null></null></null></null></null></null>	BUILDING_NO 21 44 18 <null> 12</null>	ENTRANCE <null> <null> <null> <null> <null></null></null></null></null></null>	FLOOR <null> <null> <null> <null></null></null></null></null>	· ·
POSTAL_CODE 1 1000 0 1000 0 1000 0 1000 0 1000 0	МUNICIPLITY СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ	TYPE KB. KB. KB. KB. KB. KB.	DISTRICT ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР	STREET_TYPE УЛ УЛ УЛ УЛ УЛ УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ ИВАН ДЕНКОГЛУ ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ АНГЕЛ КЪНЧЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ	STREET_NO <null> <null> 2 <null> <null></null></null></null></null>	BLOCK_NO <null> <null> <null> <null> <null> <null> <null></null></null></null></null></null></null></null>	BUILDING_NO 21 44 18 <null> 12 14</null>	ENTRANCE <null> <null> <null> <null> <null> <null></null></null></null></null></null></null>	FLOOR <null> <null> <null> <null> <null></null></null></null></null></null>	· · · · · · · · · · · · · · · · · · ·
POSTAL_CODE 1 1000 0 10	МUNICIPLITY СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ	TYPE KB. KB. KB. KB. KB. KB. KB. KB. KB. KB.	<u>DISTRICT</u> ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР	STREET_TYPE УЛ УЛ УЛ УЛ УЛ УЛ УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ ИВАН ДЕНКОГЛУ ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ АНГЕЛ КЪНЧЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ	STREET_NO <null> <null> 2 <null> <null> <null></null></null></null></null></null>	BLOCK_NO <null> <null> <null> <null> <null> <null> <null> <null> <16</null></null></null></null></null></null></null></null>	BUILDING_NO 21 44 18 <null> 12 14 12 14</null>	ENTRANCE <null> <null> <null> <null> <null> <null> <null></null></null></null></null></null></null></null>	FLOOR <null> <null> <null> <null> <null> <null></null></null></null></null></null></null>	
POSTAL_CODE 1 1000 0 1000 0 1000 0 1000 0 1000 0 1000 0 1000 0 1000 0 1000 0 1000 0	МUNICIPLITY СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ	TYPE KB. KB. KB. KB. KB. KB. KB. KB.	<u>DISTRICT</u> ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР	STREET_TYPE УЛ УЛ УЛ УЛ УЛ УЛ УЛ УЛ УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ ИВАН ДЕНКОГЛУ ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ АНГЕЛ КЪНЧЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ	STREET_NO <null> <null> 2 <null> <null> <null> <null></null></null></null></null></null></null>	BLOCK_NO <null> <null></null></null></null></null></null></null></null></null></null></null></null></null>	BUILDING_NO 21 44 18 <null> 12 14 12 14 16</null>	ENTRANCE <null> <null> <null> <null> <null> <null> <null> <null> <null></null></null></null></null></null></null></null></null></null>	FLOOR <null> <null> <null> <null> <null> <null> 1</null></null></null></null></null></null>	
POSTAL_CODE 1 1000 0 10	МUNICIPLITY СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ СРЕДЕЦ	TYPE KB. KB. KB. KB. KB. KB. KB. KB. KB.	<u>DISTRICT</u> ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР ЦЕНТЪР	STREET_TYPE УЛ УЛ УЛ УЛ УЛ УЛ УЛ УЛ УЛ УЛ	STREET ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ ИВАН ДЕНКОГЛУ ХРИСТО БЕЛЧЕВ/ГАВРИЛ ГЕНОВ АНГЕЛ КЪНЧЕВ ГРАФ ИГНАТИЕВ ГРАФ ИГНАТИЕВ	STREET_NO <null> <null> 2 <null> <null> <null> <null> <null></null></null></null></null></null></null></null>	BLOCK_NO <null> <null></null></null></null></null></null></null></null></null></null></null></null></null>	BUILDING_NO 21 44 18 <null> 12 14 16 16</null>	ENTRANCE <null> <null></null></null></null></null></null></null></null></null></null></null></null>	FLOOR <null> <null> <null> <null> <null> <null> 1 2</null></null></null></null></null></null>	

Project activities - 2. Parsing (Plan)

Parsing plan – Ataccama power in practice





Project activities - 2. Parsing (Plan components)

Parsing plan – components



SADASTRA

Project activities – 3. Cleanse & Standardize

Goal:

- Clearing different formats;
- Standardization through approximately lookups to clean up spelling errors.

			NULL	9,667	1.56%
СРДЕЦ	1	0.00%	КРАСНО СЕЛО	48,789	7.88%
СРЕДЕЦ	1	0.00%	ЛЮЛИН	48,533	7.84%
СРЕДЕД	1	0.00%	МЛАДОСТ	46,706	7.54%
СРЕДЕЦ.	1	0.00%	ТРИАДИЦА	40,732	6.58%
СРЕДЕЦИ	1	0.00%	СЛАТИНА	33,632	5.43%
СРЕДИКА	1	0.00%	ЛОЗЕНЕЦ	32,454	5.24%
СРЕЗЕЦ	1	0.00%	ВИТОША	31,671	5.11%
CPELLEL	1	0.00%	НАДЕЖДА	29,973	4.84%
СТДЕНТСКА	1	0.00%	подуяне	27,832	4.49%
СТИДЕНТСКА	1	0.00%	ИСКЪР	27,595	4.46%
СТИУДЕНТСКА	1	0.00%	ОВЧА КУПЕЛ	24,311	3.93%
СТКУДЕНТСКА	1	0.00%	СРЕДЕЦ	23,272	3.76%
СТУДЕНТСКА	1	0.00%	ВЪЗРАЖДАНЕ	23,169	3.74%
СТУДЕН	1	0.00%	ОБОРИЩЕ	22,035	3.56%
СТУДЕНКСА	1	0.00%	СЕРДИКА	21,392	3.45%
СТУДЕНТКА	1	0.00%	КРАСНА ПОЛЯНА	20,266	3.27%
СТУДЕНТСК	1	0.00%	ВРЪБНИЦА	18,295	2.95%
СТУДЕНТСКИА	1	0.00%	ИЗГРЕВ	17,985	2.90%
СТУДЕНТСКО	1	0.00%	илинден	16,536	2.67%
СТУДНТСКА	1	0.00%	СТУДЕНТСКА	14,034	2.27%
СТУДУНТСКА	1	0.00%	КРЕМИКОВЦИ	12,537	2.02%
СТЪДЕНТСКА	1	0.00%	НОВИ ИСКЪР	10,985	1.77%
Судентска	1	0.00%	ПАНЧАРЕВО	9,875	1.59%

Project activities – 3. Cleanse & Standardize (Plan)

Cleanse plan – Ataccama power in practice



Project activities – 4. Match and Merge

- Goal achieve a single view of customers
- Prerequisite defined rules for matching (matching keys)

	Cleansed data										
First	Last	G	SIN	Birth Date	Address						
John	Smith	М		1978-12-16	V3R 2A9;BC;Surrey;14618 110 Avenue						
John	Smith	М	095242434	1978-12-16	V3R 2A9;BC;Surrey;14618 110 Avenue						
John	Smith	Μ	095242434		M4X 1V5;ON;Toronto;25 Linden Street						
	Smith	М		1978-11-16							
John	Smith	Μ	095252433	1978-11-16	L3T 7M8;ON;Markham;8500 Leslie Str.						
John	Smith	Μ		1978-11-16	L3T 7M8;ON;Markham;8500 Leslie Str.						
John	Smiht		095252433	1978-11-16							
Jane	Watson	F	420347213		L3T 7M8;ON;Markham;8500 Leslie Str.						
Jane	Watson	F	420347213	1982-01-01	L3T 7M8;ON;Markham;8500 Leslie Str.						
Jane	Smith	F	420347213	1982-01-05							
J.	Smith		420347213								

SADASTRA

Project activities – 4. Match and Merge (cont.)

Cleansed data									
First	Last	G	SIN	Birth Date	Address				
John	Smith	М		1978-12-16	V3R 2A9;BC;Surrey;14618 110 Avenue				
John	Smith	М	095242434	1978-12-16	V3R 2A9;BC;Surrey;14618 110 Avenue				
John	Smith	М	095242434		M4X 1V5;ON;Toronto;25 Linden Street				

Golden record								
First	Last	G	SIN	Birth Date	Address			
John	Smith	Μ	095242434	1978-12-16	M4X 1V5;ON;Toronto;25 Linden Stree			





Project activities – 5. Enrich

		pur_zip	pur_district	std_district	titles_out_quarter	pur_quarter	
	1-	1	СТУДЕНТСКА	СТУДЕНТСКИ	KB.	ДЪРВЕНИЦА	1756
	2	1618	КРАСНО СЕЛО	КРАСНО СЕЛО	ЖК.	КРАСНО СЕЛ	1618
	3	1229	НАДЕЖДА	НАДЕЖДА	ЖК.	НАДЕЖДА III	1229
	4	7	люлин	люлин	ЖК.	люлин-іх	1373
	5	1281	НОВИ ИСКЪР	НОВИ ИСКЪР		С СЛАВОВЦИ	1281
	6	1000	ЛОЗЕНЕЦ	ЛОЗЕНЕЦ	KB.	ЛОЗЕНЕЦ	1164
Missing ZIP	7	1504	ОБОРИЩЕ	ОБОРИЩЕ	KB.	ДОКТОРСКИ ПАМЕТНИК	1504
	8	1172	ИЗГРЕВ	ИЗГРЕВ	ЖК.	ДИАНАБАД	1504
	9	1309	илинден	илинден	ЖК.	СВЕТА ТРОИЦА	1309
	10	1680	ТРИАДИЦА	триадица	ЖК.	БОКАР	1404
	11	1680	ТРИАДИЦА	триадица	ЖК.	БОКАР	1404
	12	1505	СЛАТИНА	СЛАТИНА	KB.	РЕДУТА	1505
	X	1231	НАДЕЖДА	НАДЕЖДА	ЖК.	СВОБОДА	1309
	14	1/2	ЛОЗЕНЕЦ	ЛОЗЕНЕЦ	KB.	ЛОЗЕНЕЦ	1407
	15	· 7	ЛОЗЕНЕЦ	ЛОЗЕНЕЦ	KB.	ЛОЗЕНЕЦ	1407
	16	1220	НАДЕЖДА	НАДЕЖДА	ЖК.	НАДЕЖДА II	1220
	17	,	ОВЧА КУПЕЛ	ОВЧА КУПЕЛ	KB.	ОВЧА КУПЕЛ 1	1618
	18	1618	ОВЧА КУПЕЛ	ОВЧА КУПЕЛ	KB.	ОВЧА КУПЕЛ	1618
	19	1618	ОВЧА КУПЕЛ	ОВЧА КУПЕЛ	KB.	ОВЧА КУПЕЛ	1618
	20	1592	ИСКЪР	ИСКЪР	ЖК.	ДРУЖБА-І	1592
	21	1379	ВЪЗРАЖДАНЕ	ВЪЗРАЖДАНЕ	ЖК.	СЕРДИКА	1379
	22	1505	ОБОРИЩЕ	ОБОРИЩЕ	KB.	ПОДУЯНЕ-ЦЕНТЪР	1320
	23	;	ИСКЪР	ИСКЪР	ЖК.	ДРУЖБА-Г	1330
	24	1618	КРАСНО СЕЛО	КРАСНО СЕЛО	ЖК.	БЪКСТОН	1618
	25	1618	КРАСНО СЕЛО	КРАСНО СЕЛО	ЖК.	БЪКСТОН	1618
	26	;	илинден	илинден	KB.	ЗАХАРНА ФАБРИКА	1618
	27	,	ВИТОША	ВИТОША	ЖК.	БЪКСТОН	1618
	28	;	КРАСНО СЕЛО	КРАСНО СЕЛО	ЖК.	БОРОВО	1680



Project activities - 5. Enrich (Plan)

Goal - elaborate with additional information from reference sources



≪ADASTR

31

Project Delivery Good Data = Good Business

- Cleansed Data
 - Corrected typing errors, removed dummy characters, etc.
- Merged Data
 - No duplications of records
- Elaborated Data
 - All addresses are completed with ZIP codes
- Refined Data Quality Rules
 - Implementing rules to observe data quality
- Data Quality Report
 - Current Data Quality status and quantified DQ issues
 - Next steps of DQ improvement



Key Takeaways

Good company

 Adastra is a solid world company and gives a career opportunity for motivated young people without experience

Good perspectives

Data Quality is an increasing niche and perspective market

Good skills

 Gain both business knowledge and technical skills with best-ofbreed technologies



Thank You



ADASTRA Bulgaria

29 Panayot Volov Str. 1527 Sofia, Bulgaria Tel: +359 2 960 00 30

2 Dunav Str. 9000 Varna, Bulgaria Tel: +359 2 960 29 95

www.bg.adastragrp.com

infobg@adastragrp.com

jobsbg@adastragrp.com



ADASTRA GROUP North America 8500 Leslie Str., Suite 600 Markham, Ontario CANADA L3T 7M8 Tel: +1 905 881 7946 info@adastragrp.com

