Dhan Patel :- Senior Product Manager Information Governance :- IBM Software Group

# The Essentials of Data Discovery: Do You Know Where Your Data Is?





#### **Important Disclaimer**

for a stharter planet C

THE INFORMATION CONTAINED IN THIS PRESENTATION IS PROVIDED FOR INFORMATIONAL PURPOSES ONLY.

WHILE EFFORTS WERE MADE TO VERIFY THE COMPLETENESS AND ACCURACY OF THE INFORMATION CONTAINED IN THIS PRESENTATION, IT IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED.

IN ADDITION, THIS INFORMATION IS BASED ON IBM'S CURRENT PRODUCT PLANS AND STRATEGY, WHICH ARE SUBJECT TO CHANGE BY IBM WITHOUT NOTICE.

IBM SHALL NOT BE RESPONSIBLE FOR ANY DAMAGES ARISING OUT OF THE USE OF, OR OTHERWISE RELATED TO, THIS PRESENTATION OR ANY OTHER DOCUMENTATION.

NOTHING CONTAINED IN THIS PRESENTATION IS INTENDED TO, OR SHALL HAVE THE EFFECT OF:

- CREATING ANY WARRANTY OR REPRESENTATION FROM IBM (OR ITS AFFILIATES OR ITS OR THEIR SUPPLIERS AND/OR LICENSORS); OR
- ALTERING THE TERMS AND CONDITIONS OF THE APPLICABLE LICENSE AGREEMENT GOVERNING THE USE OF IBM SOFTWARE.



#### Agenda

- Overview
- Discovery Deep Dive
  - -Referential Integrity & Business Object discovery
  - -Transformation discovery
  - -Sensitive/Critical Data discovery
  - -Discovery for Data Consolidation
- Customer Case Studies

for a sharter planet



#### Overview





#### You can't ----- what you don't understand



- Increasingly distributed
- Complex, poorly documented data & relationships within & across sources
- Data not understood because:
  - Corporate memory is poor
  - Documentation is poor or nonexistent
  - Logical relationships (enforced through application logic or business rules) are *hidden*



#### Data Discovery: An Invaluable Data Analysis Tool



- Automated analysis of data and relationships for *complete understanding* of data assets:
  - Characterizes data elements within a Source
  - Identifies *relationships* that link data elements into "business entities" *within* a source
    - Customer, counterparty, invoice
  - Identifies complex logic that relates multiple sources



#### Poor Understanding = Unpredictable Project Deployment



#### Data Discovery:

#### Automation that accelerates time to value:

a strarter

#### Data Growth Management:

Automates discovery of referential integrity and business objects

#### Data Consolidation, Integration & Migration :

Discovers transformation and business logic between data sources

Prototypes empty targets from the combination of many data sources

#### Data Privacy:

Discovers hidden sensitive data

#### What is unique?

- Analyzes data values and patterns and produces actionable results
- Discovers complex relationships within and between data sources
- Patented approach with greatest level of automation in the industry



for a marter planet M



**Discovery Deep Dive** 



## Information Management for a sharter planet



#### **Referential Integrity & Business Object Discovery**





Referential Integrity & Business Object Discovery: From knowing nothing to knowing what you need to know



#### se smarter planet 2 c n 1 s marter

Discover Complete Business Object for Optim Projects: Archiving, Test Data, Application Retirement (Optim)





#### **Business Object Volume Analysis**

- Prototyping archive strategy
  - Define a selection condition
  - Identify tables that will be kept on source or archived away.
  - How much volume do I archive away?





#### The Crown Jewels: Transformation Discovery



for a 🗱 narter planet 🛄



#### Transformation Discovery By Comparing Data: Aligns Rows Based on Correlations

a stharter planet

Step 1: Discovery Engine analyzes the data values to automatically discover the columns that align rows across data sources:





#### Transformation Discovery By Comparing Data : Column Mapping

a anar

Step 2: With rows now aligned, analyzes the data values to automatically discover complex business rules and transformations:										
Case:If age<18 and Sex=M then 0 If age<18 and Sex=F then 1If age>=18 and Sex=M then 2 If age>=18 and Sex=F then 3 $=$ Demo1										
Table 1							Table 25			
Row	Member	SS #	Age	Phone	Sex		ID	Demo1		
1	595846226	123-45-6789	15	(123) 456-7890	М		595846226	0		
2	567472596	138-27-1604	8	(138) 271-6037	F		567472596	1		
3	540450091	154-86-4196	22	(154) 864-1961	М		540450091	2		
4	514714372	173-44-7900	55	(173) 447-8996	F		514714372	3		
5	490204164	194-26-1648	4	(194) 261-6476	F		490204164	1		
6	466861109	217-57-3046	66	(217) 573-0453	М		466861109	0		
•	•	•	•	•	•		٠	•		
•	•	•	•	•	•		•	•		
987,623	444629628	243-68-1812	25	(243) 681-8107	F		444629628	3		
987,624	423456789	272-92-3629	87	(272) 923-6280	M		423456789	2		
		1						1		

## Transformation Discovery By Comparing Data : Measuring Mapping Rule Strength

a sinar

		to auto	mate disco	overy of <b>unknown</b> oves QA early i	data in n the p	consistencies process		
Case:	lf age<18 If age<18	and Sex=M and Sex=F	then 0 then 1	If age>= If age>=	<mark>18 and</mark> 18 and	d Sex=M then 2 d Sex=F then 3	= Der	no1
Table 1							Table 25	
Row	Member	SS #	Age	Phone	Sex	]	ID	Demo1
1	595846226	123-45-6789	15	(123) 456-7890	м		595846226	0
2	567472596	138-27-1604	8	(138) 271-6037	F		567472596	1
3	540450091	154-86-4196	22	(154) 864-1961	М		540450091	2
4	514714372	173-44-7900	55	(173) 447-8996	F		514714372	3
5	490204164	194-26-1648	4	(194) 261-6476	F		490204164	1
6	466861109	217-57-3046	66	(217) 573-0453	м		466861109	0
•	•	•	•	•	•		•	•
	•	•		•			•	•
987,623	444629628	243-68-1812	25	(243) 681-8107	F		444629628	3
•	•	272 02 2620	87	(272) 023-6280	M		422456780	

#### Transformation Discovery By Comparing Data: Discoverable **Transformations**

- Scalar
- One to one
- Substring
- Concatenation

- Constants
- Tokens

- Conditional Logic
  - Case Statements
  - Equality/Inequality
  - Null Conditions
- In/Not In
  - Conjunctions

- Joins
- Inner
- Left Outer
- Average
  - Minimum

• Sum

Maximum

Aggregation

- Column Arithmetic
- Add
- Subtract
- Multiply
- Divide

- **Reverse Pivot**
- **Cross-**Reference
- **Custom Data** ٠ **Rules**

	Mapping Studio - LOCALHOST - HR Master Data	Management			
	Project View Tools Map Help				
	Home 🥝 Data Sets 🗳 Column Analysis 🧉	🍦 PF Keys 🗳 Data Objects 🗳 Tar	rget Matches 🗳	Maps	
	Maps: DOMap_DOMap1_DO_HQ_EMP	_to_DO_W_EMPS1_W_EMPS		HR	Master Data Management   No Activi
	Summary O Joins 😜 Bindings 🔿 ۱	Where Clause Gransformations O	Reverse Pivots		
	🕂 💥 🔯   🥅 Show Data 🔹 💽 Export Hits 🔹   📢 Fil	ter Matches 🛛 🎼 Generate Data Rule 🛛 🌌			
CASE WHEN HQ_EMP.STATUS in ( ' Current', 'Fired', 'Resigned') THEN HQ EMP.TERMINATION DATE ELSE	Target Column   Imaget Column	Primary Transformation   destudieR_EMPCOVEC_D_0_H01   HQ_EMP_ITTLE_OF_CORTESY   HQ_EMP_ITTLE_OF_CORTESY   HQ_EMP_ITTLE_OF_CORTESY   HQ_EMP_INAMEI   HQ_EMP_INAME2   HQ_EMP_INAME2   HQ_EMP_INAME2   HQ_EMP_INAME2   HQ_EMP_INAME2   HQ_EMP_INAME2   CASE_WHEIN_QEMP_STATUS in (' 1   HQ_EMP_IO_D0H   CASE_WHEIN_QEMP_STATUS in (' 1   CASE_WHEIN_QEMP_STATUS in (' 1   CASE_WHEIN_QEMP_STATUS in (' 1	HILS 000.00 % (20/20) 90.00 % (20/20) 100.00 % (20/20) 100.00 % (20/20) 100.00 % (20/20) 100.00 % (20/20) 100.00 % (20/20) 100.00 % (20/20) HQ_EMP.TERMINAT	Misses 0.03% (0/20)   0.00% (0/20) 0.00% (0/20)   0.00% (0/20) 0.00% (0/20)   0.00% (0/20) 0.00% (0/20)   0.00% (0/20) 0.00% (0/20)   0.00% (0/20) 0.00% (0/20)	Columns   [All Columns]   Source Columns   MAME1   NAME2   SN   TITLE   Do6   DO6   DO6   DO6   StATUS
HQ_EMP.RETURN_DATE END					Columns Data Rules
	Validate Step Approve Step				Re-Run Step Run Next Steps.
	Maps: DOMap_DOMap1_D0_HQ_EMP_to_D0_W_EMPS1_W	EMPS			



#### **Discovery for ETL and Lineage Metadata Initiatives**





#### **Transformation Discovery Example**

a Sanar

			Transf	ormation S	Stats		
DO FEED2 To DO PRODUCT	DATA_1566011	1033			Pro	i   <u>No Ac</u>	tivity
Summary 🔿 Joins 🛆 Bindings (	→ Where Clause △	Transformations O Revers	e Pivots				
•	•						
Q • 9 🗿 🛍 -   🕂 💥	🙀   🗙 Show Misses 👻	👔 Export <u>H</u> its 🕞 🕸 Generate Data Rul	e				
# Target Column	Primary Transformatio	n	Δ	Hits	Misses	Notes	
12 CPN_FLOAT_RATE_RESET_FRE	CASE WHEN FEED2.COUPO	N_TYPE in ('FLOAT', 'VARIABLE') THEN '	Monthly' ELSE n	96.49 % (660/684)	3.51 % (24/684)	2	
COUPON_LAST_REGULAR_PMT	CASE WHEN FEED2.SECURI	TY_ALIAS in ( 107450, 141458, 14463	1, 159719, 15	99.42 % (680/684)	0.58 % (4/684)		
+ PRODUCT_IDENTIFIER	datarule(DR_PRODUCT_IDE	NTIFIER_0, FEED2.SECURITY_DESCR)		100.00 % (684/684)	0.00 % (0/684)		
+ 14 ABG_INDUSTRY_NAME	FEED2.BA_INDUSTRY_SECT	OR		98.39 % (673/684)	1.61 % (11/684)		
3 COUPON_CURRENCY	FEED2.COUPON_CURRENC	Y		98.39 % (673/684)	1.61 % (11/684)		
	· · · · · · · · · · · · · · · · · · ·	1		98.25 % (672/684)	1.75 % (12/684)		≡
Discovered cor	npiex	pr .		88.45 % (605/684)	11.55 % (79/684)		
transformation (Case	Statement)			98.25 % (672/684)	1.75 % (12/684)		
	Statement	VTION		59.21 % (405/684)	40.79 % (279/684)		
				98.39 % (673/684)	1.61 % (11/684)		
+ 16 ISSUE_PRICE	PX			98.25 % (672/684)	1.75 % (12/684)		
21 ORIGINAL_AMOUNT	AMT			98.39 % (673/684)	1.61 % (11/684)		
22 OUTSTANDING_AMOUNT	FEED ANDING_AMT			100.00 % (684/684)	0.00 % (0/684)	2	
23 COMPUTED_OUTSTANDING_A	FEED2.0 TANDING_AMT	* 2		37.72 % (258/684)	62.28 % (426/684)		
+ 17 PAR_VALUE_PER_SHARE	FEED2.UNIT_PAR_AMT * 0.	1		98.25 % (672/684)	1.75 % (12/684)		~
CASE NHEN FEEDS COUDON TYPE in	( LELOATI IVART	NRIEL THEN Monthlul FI	F pull END		•		
CASE WHEN FEED2.COUPON_TIPE IN	( FLOAT , VARIA	RDLE , THEN MONCHTY' ELS	SE HUII END				



#### Transformation Discovery Example



a 🚌 narter planet 🗰



#### Pattern Based Sensitive Data Discovery Example: SSN

#### InfoSphere Discovery Classified Columns View

Ŀ	- MC		Master	Card	Yes		Yes		No			
E	NE	A	Notes E	mail Address	No		Yes		No			
E	- NI	=	Spanish	NIF	No		Yes		No			
E	NII	VO	UK NIN	0	No		Yes		No			
E	PA	N	Passpo	rt Number	No		Yes		No			
E	PN		Person	al Name	No		Yes		No			
Ŀ	- SI	J.	Canadi	an SIN	No		Yes		No			
	- SS	N	US Soci	al Security Number	Yes		Yes		No			
	-	Assigned Excluded										
		Column Metadata					Column Cla	ssificatio	חו	Statistics		
	0	Name		Table		Data Set	Hit Rate	Ap	Method	Cardinality	Selectivity	Null Count
	Þ	COMMENTS		ALL_POSITIONS_F		Data Set 1	50%		Discovered	4	0.00	3641
		MRC_ACCTD_AMOUNT		MK_ADJTMNT		Data Set 1	35.2941%		Discovered	12	0.04	257
		MRC_EXCHANGE_RATE		MK_KE_TRX_70		Data Set 1	52.5%		Discovered	45	0.00	9397
		MRC_EXCHANGE_RATE		MK_SCHED		Data Set 1	33.2075%		Discovered	47	0.00	9590
		MRC_ACCTD_EARNED_DISC_	TAKEN	MK_SHOU_APPS_A	۱LL	Data Set 1	100%		Discovered	1	0.00	200
		MRC_ACCTD_UNEARNED_DIS	5C_T	MK_SHOU_APPS_A	۱LL	Data Set 1	100%		Discovered	1	0.00	200
		ID_PIN		NK_JE_HDRS		Data Set 1	100%		Discovered	23	0.00	0
		BILLING_ID		PRFL_3_CUST		Data Set 1	39.777%		Discovered	269	1.00	0
E		L	URL		No		Yes		No			
E	US	PHN	US Pho	ne Number	Yes		Yes		No			
E	- US	SC	US Stat	e Code	No		Yes		No			









#### Beyond "Scanning for Pattern": additional sensitive data discovery

### Discovery finds other types of sensitive data

- Sensitive data that do not have distinctive pattern
- Sensitive data that are partial or hidden (malicious)
- Sensitive data embedded in free text

#### How do we find them?

- data similarity exact value matching
- data similarity -- fuzzy value matching
- metadata similarity including known classification
- data relationships



#### Scan for pattern: custom "sensitive", custom "algorithm"

### Define "custom sensitive" data in Discovery

## • How do we find custom sensitive data?

- User defined algorithm -- Once deployed, behave the same as "built-in"
  - Algorithm is run as part of profiling
  - Hit/miss metrics and data view will be available on custom sensitive.
- data similarity exact value matching
- data similarity -- fuzzy value matching
- metadata similarity including known classification
- data relationships



From "Sensitive Data" Discovery" to "Critical Data" Discovery: Using Discovery to connect Business to IT

- Discover "critical data elements" (CDEs) and map them to Business Glossary
  - Data driven term mapping
- Once we find these CDEs, explore its surroundings to identify other data relevant to the business of interest.

#### Sensitive/Critical Data Discovery For Compliance and Test Data Management Projects









Masking privacy data in test data extract



### Complete methodology for Consolidation Prototyping

strater blar





#### Data Inventory: Cross System Overlap

🗈 Discovery Studio - LOCALHOST - HR Data Overlap													
Project \	View Tool	s Help											1
Home (	🔿 Data S	ets 🔿 Column		PF Ke	vs 🔿 Dat	a Objects	📿 Overlan	5					
(	<i>J b a a b</i>					000,000		-					
Overlaps													
EMP_ALL													
Column Sun	nmary Data	a Set Overlaps								ack 🔻 Show All	Items	-	
0		- 🗷 🏟 - 👧											
~		. 5. 6. 4											
Summary	Table	Calum	Column Number	CDE	Cardinality	Coloctivity	New Nulls Of	Value Overlap	with Benefits	Value Overlap v	with Payroll	Value Overlap wi	th Skillset
	EMP ALL	EMP ID	Coldmin Number		Cardinality 50	1.00	100 %	EMP_ALL 82.%	100 %	EMP_ALL 82 %	100 %	24 %	100 %L
EMP ALL	EMP ALL	DEPT ID	2		48	0.96	100 %	0 %	0 %	75 %	90 %	25 %	100 %
EMP_ALL	EMP_ALL	TITLE	3		20	0.40	100 %	0 %	0 %	75 %	75 %	50 %	100 %
EMP_ALL	EMP_ALL	MANAGER ID	4		6	0.12	100 %	83 %	<u>12 %</u>	83 %	12 %	83.%	100 %
EMP_ALL	EMP_ALL	FIRST NAME	5		<u>49</u>	0.98	100 %	Identifie	s overlar	s across	data s	ources	100 %
EMP_ALL	EMP_ALL	LAST NAME	6		<u>50</u>	1.00	100 %	<u>0 %</u>	<u>0 %</u>	<u>70 %</u>	<u>85 %</u>	<u>24 %</u>	100 %
EMP_ALL	EMP_ALL	ADDRESS 1	7		<u>50</u>	1.00	100=%	Allows a	nalyst to	select c	ritical d	lata eleme	ents (CDE
EMP_ALL	EMP_ALL	ADDRESS 2	8		1	0.02	100 %	for future	<u>**</u> مورر د	<u>0 %</u>	<u>0 %</u>	<u>0 %</u>	<u>0 %</u>
EMP_ALL	EMP_ALL	ADDRESS CITY	9		<u>50</u>	1.00	100 %		0 0 <u>0 %</u>	<u>82 %</u>	<u>100 %</u>	<u>24 %</u>	<u>100 %</u>
EMP_ALL	EMP_ALL	ADDRESS STATE	10		<u>29</u>	0.58	100_%	Determi	ne what	attributes	are si	ipersets a	nd subset
EMP_ALL	EMP_ALL	ADDRESS COUNTRY	11		<u>44</u>	0.88	100 %	botwoon		84 %	<u>100 %</u>	27%	100 %
EMP_ALL	EMP_ALL	ADDRESS ZIP	12	~	<u>50</u>	1.00	100 %	Dermeer	sources	<u>82 %</u>	<u>100 %</u>	<u>24 %</u>	100 %
EMP_ALL	EMP_ALL	PHONE HOME	13		<u>50</u>	1.00	100 %	0%	0.%	82 %	<u>100 %</u>	24.%	<u>100 %</u>
EMP_ALL	EMP_ALL	PHONE CELL	14		<u>50</u>	1.00	100 %	0%	0%	0%	0%	24 %	100 %
EMP_ALL	EMP_ALL	PHONE WORK	15		<u>50</u>	1.00	100 %	0%	0%	0%	0%	24 %	100 %
EMP_ALL	EMP_ALL		16		<u>48</u>	0.96	100 %	0%	0%	0%	0%	23 %	100 %
EMP_ALL	EMP_ALL	EC PHONE 1	17		50	1.00	100 %	0 %	100.00	0 %	100.01	24 %	100 %
EMP_ALL		DATE HIKED	18		50	1.00	100 %	82 %	100 %	<u>82 %</u>	100 %	42.94	100 %
EMP_ALL_		BILL BY	19		<u>23</u>	0.46	100 9/20		11.%	<u>03 %</u>	<u>95 %</u>	<u>43 %</u>	100 %

#### Data Inventory: Critical Data Element Alignment Across sources

0	Target Table S	chema 🔵 Source Map	oping 🖕 Vi	nified Col	lumn	Analysis	O Match a	nd Merge A	nalysis			
	$\sim$											
То	Fotal Source Row Count: 189											
٩	🖌 🔹 🖓 🚰 🏚 🖌 🛅 Preview Data 📑 Value Frequency 👻 🥅 All Mappings 🔂 Export Table											
Tar	arget Table Statistics											
#	Name			D	)ata Ty	/pe	Data Sets	Cardinality	Selectivity	Min	Max	
+	1 F_NAME			Va	archar	(30)	100.00 % (3/3)	108	0.5714285	ABBY	YVONNE	
+	2 L_NAME			Vä	archar	·(30)	100.00 % (3/3)	89	0.4708994	AARON	WELCH	
Ξ	3 SOCIAL			Va	archar	·(15)	100.00 % (3/3)	145	0.7671957	103-58-1068	396694138	
	Mapping		Statistics									
	Data Set	Transformation	Cardinality	Selectivity	Ν	Min	Max	Complete				
	Target Table		145	0.7671957	76 1	103-58-1068	396694138					
	Community	COMMUNITY_BRCH.SOCIAL	52		1 1	103-58-1068	387-29-3234	100.00 % (	5			
	Region	REGION_BRCH.SSN	51		1 1	106619513	396694138	100.00 % (	5			
	CRM	CRM_BRCH_1A.TAX_ID	86		1 1	103-58-1068	396-69-4138	100.00 % (8	8			
+	4 CITY			va	archar	(20)	100.00 % (3/3)	96	0.5079365	ADONA	YEADON	
+	5 STATE_PROV			Va	archar	(12)	100.00 % (3/3)	42	0.2222222	AL	WV	

- Provides profiling data for the union of all of the sources mapped to the virtual target
- Shows profiling results for individual source attribute and for the union of all sources

#### Unified Modeling: Bottoms Up Modeling of Target Schema



#### Sources to Target Mapping Discovery: using both data and metadata

Q Sour	ce Mapping - Sug	gest Transforma	tions: NEWTAB	LEO: R	egion				
9	• 💡 🛔	🏻 🎦 🕶							
Target Co	olumns	Source Tables					Statistics		
Select	Target Column	Table	Column	CDE	Cla	Selectivity	Name Match	Row Hit Rate	Value Hit Rate
	F_NAME	REGION_BRCH	MI			0.96		<u>7.84 % (4/51)</u>	<u>6.12 % (3/49)</u>
	F_NAME	REGION_BRCH	FN			0.96		52.94 % (27/51)	51.02 % (25/49)
	L_NAME	REGION_BRCH	LN			0.86		72.55 % (37/51)	<u>68.18 % (30/44)</u>
	CITY	REGION_BRCH	CITY			0.96	✓	72.55 % (37/51)	71.43 % (35/49)
	STATE_PROV	REGION_BRCH	STATE			0,51		<u>88.24 % (45/51)</u>	80.65 % (25/31)
Record 1	of 5 🖣								•
								Бк	Help
Name matches are also shown Overlap percentages between a data source to the virtual target									

- Map sources to the target schema based on name and value matches (from the overlap analysis)
- Each source mapped adds more data to the "virtual target" and enables more value matching from the overlap analysis

#### Match and Merge Analysis

Data Set	F_NAME	L_NAME	SOCIAL	CITY	STATE_PROV
<merged></merged>	SERGIO	ELLIOTT	230-24-7823	PARKERS PRAIRIE	GA
COMMUNITY	SERGIO	ELLIOTT	230-24-7823	PARKERS PRAIRIE	GA
COMMUNITY	ELEANOR	ELLIOTT	295-78-2505	LOOMIS	PA
COMMUNITY	MARSHAL	ELLIOTT	272-72-7078	NORTH OAKS	CA
COMMUNITY	ROGER	ELLIOTT	218-28-1759	ORCHARD HEIGHTS	NV
CRM	SERGIO	ELLIOTT	230-24-7823	PARKERS PRAIRIE	GA
CRM	ELEANOR	ELLIOTT	295-78-2505	LOOMIS	PA
CRM	MARSHAL	ELLIOTT	272-72-7078	NORTH OAKS	CA
CRM	ROGER	ELLIOTT	218-28-1759	ORCHARD HEIGHTS	NV
REGION	SERGIO	ELLIOTT	230-24-7823	PARKERS PRAIRIE	GA
REGION	ELEANOR	ELLIOTT	295-78-2505	LOOMIS	PA
REGION	MARSHAL	ELLIOTT	272-72-7078	NORTH OAKS	CA
REGION	ROGER	ELLIOTT	218-28-1759	ORCHARD HEIGHTS	NV

Data Set	F_NAME	L_NAME	SOCIAL	ατιγ	STATE_PROV
<merged></merged>	SERGIO	ELLIOTT	230-24-7823	PARKERS PRAIRIE	GA
COMMUNITY	SERGIO	ELLIOTT	230-24-7823	PARKERS PRAIRIE	GA
CRM	SERGIO	ELLIOTT	230-24-7823	PARKERS PRAIRIE	GA
REGION	SERGIO	ELLIOTT	230-24-7823	PARKERS PRAIRIE	GA 🔷

 Match and Merge analysis allows user to quickly prototype using different columns as the matching key across sources

- Example of overmatching when using a match key of L\_Name (conflicts in red)
  - Example of matching when match key being prototyped is Social Security Number





#### **Customer Case Studies**

for a smarter planet M





#### Case Study: International Truck Manufacturer

#### **Business Challenge:**

- Improve customer experience
- Reduce number of financial systems
- Address data quality once for multiple large migration projects

#### Approach:

Consolidate various systems into a single master

#### **Technical Issues:**

- 14 customer, 9 vendor and 6 materials data sources
- Manual data analysis put project at risk
  - Required 6 months to cross-analyze just 2 systems
  - Assumptions were often incorrect and documentation was spotty
  - Each analyzed attribute involved 3-7 business and data analysts
  - Manual combination of cross system queries was error prone
  - QA process for cross-system analysis was almost non-existent





#### Case Study Results: Using InfoSphere Discovery for Cross-System Analysis

- 12 months of work completed in only 1 month
  - -~\$180K savings in first month alone using Exeros
  - Increased project predictability and lowered risk
- Assumptions about business rules and relationships are now tested and validated against millions of rows before decisions are made
- One person can answer questions in an hour that used to take three or more people several days to answer
- Less SME time is now required and less time is spent coordinating discussions between different SMEs



#### IBM

#### KeyBank

" After the product was purchased, the Key Bank team targeted applications requiring source-to-target mapping from the source systems to the enterprise data warehouse. The previous approach, requiring crafted queries and manual effort, was expected to take nine staff months. With the cross-system data analysis workbench, Exeros (now InfoSphere) Discovery, the team was able to do the mapping automatically in two weeks – a significant reduction in cost and effort but, more importantly, a reduction in time to business value.

> Source: Data Governance From Policy to Practice BeyeNETWORK Research Report Case Study: **Key Bank** David Loshin and BeyeNetwork 2008



Case Study: Data Lineage and Traceability

**Chief Data Architect** 

Regional Commercial Bank

A lineage mapping project as part of our MDM program was **estimated to cost \$120K** and take three months, **was completed for only \$10K** in one week because of Discovery software."



#### Case Study: IT Asset Master

- Problem:
- Couldn't answer simple questions:
  - Who owns an IT asset?
  - How much does it cost us to operate our current Infrastructure?
  - How much are we underutilizing our assets?
- Project:
- Consolidate 24 asset management systems into a single asset master
- Manual results:
- 7,520 person hours at \$110/hour (~\$825K) spent manually mapping
- Volume of data and complexity was too great
- Still couldn't answer the questions
- InfoSphere Discovery results:
- 2 weeks of elapsed time

Vice President Charlotte, NC based Commercial bank

One person using Discovery software, who knew nothing about our environment accomplished more in a two week pilot than 9 subject matter experts were able to complete over 9 months."

#### IBM InfoSphere<sup>™</sup> Discovery – Invaluable Data Analysis Tool

- Accelerate deployment of your information agenda projects:
  - Improves accuracy, predictability and repeatability
  - -Speeds project data analysis by as much as 10 times
  - -Minimize SME Time



Information Management





for a 🚎 arter planet 🗰