



Data Quality II: Tools

Yves Bontemps
SmaIS-MvM, Section Recherches
21 septembre 2006



- Introduction
- Outils de Data Quality: concepts
 - Data Profiling
 - Standardisation
 - Matching
 - Monitoring
- Outils commerciaux
 - Case study
 - Architecture
- Conclusion

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

2



Introduction



21/09/2006

Ring!
Ring!
Ring!



Data Quality II: Tools
Y. Bontemps - Recherches

3



Introduction



21/09/2006

Allo?



Data Quality II: Tools
Y. Bontemps - Recherches

4



Introduction

Bonjour, M. Bontemps. Désirez-vous diminuer votre facture d'électricité?

Oui, bien sûr.

Devenez client chez GreenVolt!

Je suis déjà client chez vous!

Oups. Pardon.

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

5



Introduction

! ⚡ \$ @ !

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

6



Introduction

FACTURE
M. Yves Bontemps
de la Loi, 9/2
7100 Haine St-Paul



21/09/2006

M. Yves Bontemps
Rue de la Loi, 9 bte 2
7100 La Louvière



Bottin

Data Quality II: Tools
Y. Bontemps - Recherches

7



Qualité des données Définition

Fitness for use

- Remarques:
 - Fitness vs Perfection
 - Coûts-bénéfices
 - Use → présent & futur (!)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

8



Introduction

Impacts/enjeux

- National Firearms Licensing Management Systems (UK)
 - *"During the pilot there were a number of data quality issues, which meant the system was returning errors, so the system was declined"*
 - *"If the Home Office really is incapable, over a period of eight years, of computerizing something as straightforward as a few hundred thousand firearms records, then it does suggest that they do not have a hope of making a success of the introduction of the national identity card scheme"*
- <http://www.computing.co.uk/computing/news/2148820/delay-gun-register>

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

9



Introduction

Impacts/enjeux

- Voter's registration system in California
 - Registration system of all voters, based on identification (driver's license). Checked against Calif. Dept of Motor Vehicles database.
 - *"The rigorous system will reject applications whose data doesn't exactly match the confirming documents. Even small discrepancies, such as a missing middle initial, could cause an application to be rejected."*
 - *"The voter database has "been a disaster for anyone who is trying to register for the first time or reregister because they moved, got married and need to change their name or change parties,"*
- <http://www.computerworld.com/article/0,9891,110353,00.html>

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

10



Introduction

Impacts/enjeux

- Criminal Records Bureau (UK)
 - Check that someone has no criminal record prior to appointment (esp. unsupervised contacts with children).
 - *"The Criminal Records Bureau's first and foremost priority is to help protect children and vulnerable adults"*
 - *"The Criminal Records Bureau is only as effective as the information it can access."*
 - *Liberal Democrat home affairs spokesman Nick Clegg said the errors took "Home Office incompetence to new absurd levels". He added: "This latest fiasco will erase the last bit of public confidence in the Home Office."*
 - <http://news.bbc.co.uk/1/hi/uk/5001624.stm>

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

11



Introduction

Impacts de la qualité

- Coûts de correction (usine fantôme)
- Risques accrus → nouveaux dével.
- Décisions erronées
- Perte de confiance
- Abandon/Rejet du système

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

12



Introduction

Dimensions de la qualité

- Dimensions
 - Pertinente (*Relevant*)
 - Précise (*Accurate*)
 - Fraîche (*Timely*)
 - Complète (*Complete*)
 - Comprise (*Understood*)
 - Digne de confiance (*Trustworthy*)
- Contexte = utilisation des données

21/09/2006

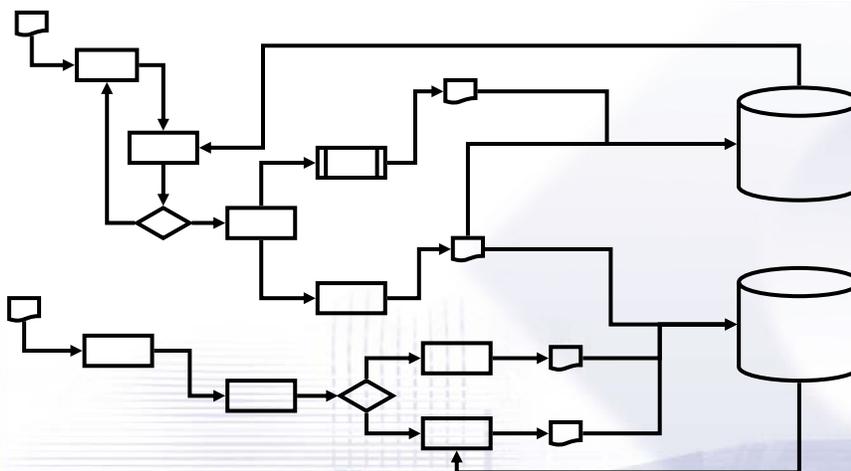
Data Quality II: Tools
Y. Bontemps - Recherches

13



Introduction

Production des données



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

14



Introduction

Facettes d'une solution

Management

- Exec. DQ group
- Data Mgmt
- Data stewards
- ...

Information

- Meta-data
- Data Models
- Requirements
- Metrics, ...

Data as an Asset

Activités

- Root-cause analysis
- Data Integration
- Software Eng.
- Business Proc. Eng.
- Cleansing, etc.

21/09/2006

Bénéfices

- Utilisateurs
- Développeurs
- Clients
- Affaires
- ...

Data Quality II: Tools
Y. Bontemps - Recherches

15



Introduction

Outils Data Quality

- Marché existant et en forte croissance
- Question
 - Fonctionnalités proposées ?
 - Coûts vs Bénéfices ?
 - Par rapport dével. in-house?
- Place dans une approche globale?

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

16



- Introduction
- **Outils de Data Quality: concepts**
 - Data Profiling
 - Standardisation
 - Matching
 - Monitoring
- Outils commerciaux
 - Case study
 - Architecture
- Conclusion

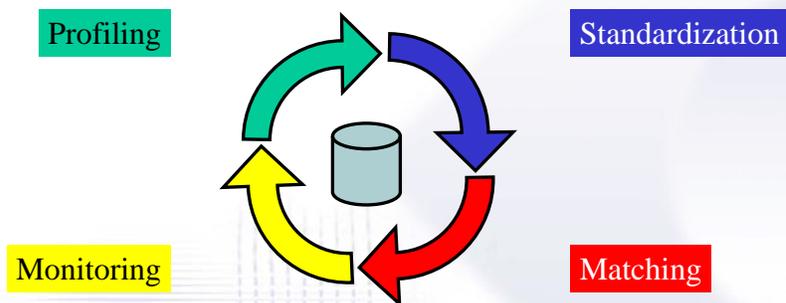
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

17



Cycle d'utilisation



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

18



Cycle d'utilisation

Profiling



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

19



Data Profiling



- Expertiser l'état d'une base de données
 - Structure
 - Contenu
 - Validité < *Accuracy*
- Bird's eye view/ Rapport résumé



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

20



Data Profiling



- Information **sur** les données:
 - 10% des champs NUM_TEL sont vides
 - Le domaine de HEURES_REF est 3600 à 4000.
- **Pas** information **à partir** des données:
 - Les ouvriers de la construction sont significativement surreprésentés dans les accidents du travail.
 - 10% des contrats de travail dans l'enseignement durent moins de 3 mois.

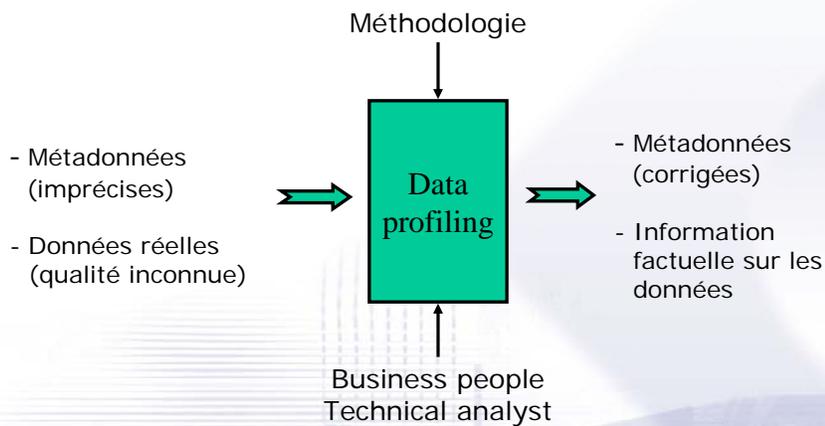
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

21



Data Profiling Modèle



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

22



Data Profiling

Problèmes



- Métadonnées:
 - Inexistantes
 - Imprécises
 - Plus à jour
- Données:
 - Incohérentes avec métadonnées
 - Invalides/incorrectes
 - Contenu différent des attentes

21/09/2006

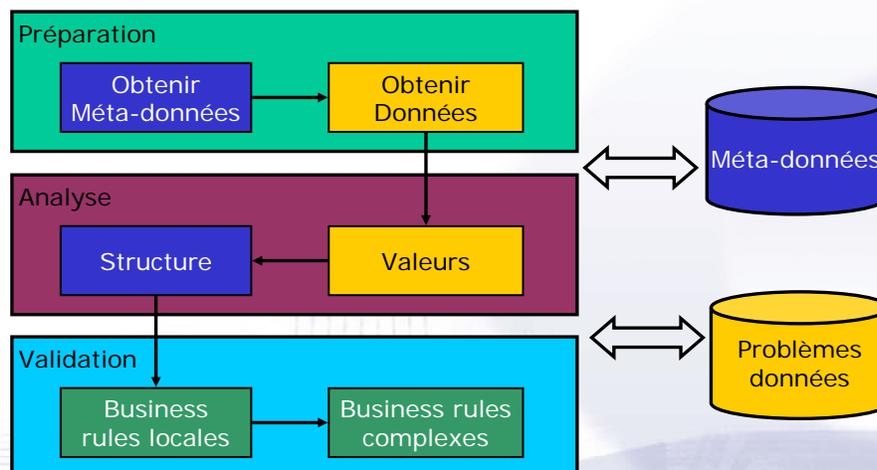
Data Quality II: Tools
Y. Bontemps - Recherches

23



Data Profiling

Processus



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

24



Data Profiling Notes



- Préparation
 - Rassembler les métadonnées, les participants, etc.
 - Extraction des données (!)
 - Effort principal
- Analyse & Validation
 - 2 phases pour chaque étape:
 - Discovery
 - Assertion testing

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

25



Data Profiling Préparation (Méta-données)



- Systèmes de gestion de l'information
 - Dictionnaires des données (ex: glossaires)
 - Meta-data repositories
- Data definitions
 - Copybooks COBOL
 - Catalogs
- Logiques/Règles business
 - Programmes
 - Analyses fonctionnelles
 - Instructions aux utilisateurs
- People
 - DBA, Data Architect, Business Analyst

21/09/2006

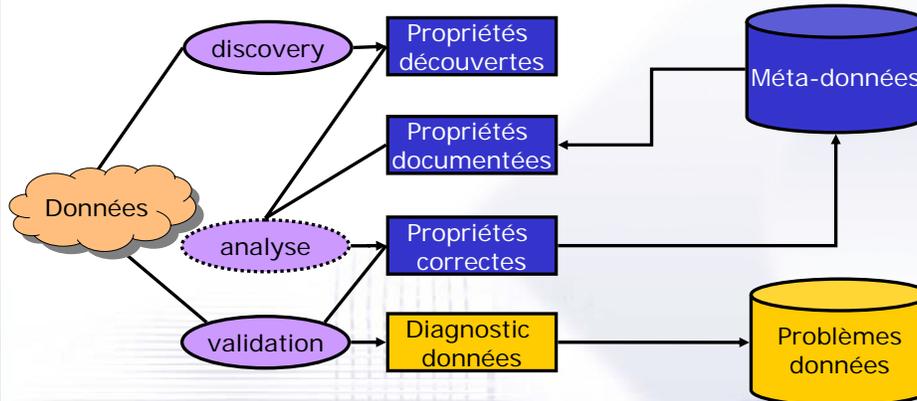
Data Quality II: Tools
Y. Bontemps - Recherches

26



Data Profiling

2 phases: exemple (valeurs)



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

27



Data Profiling

Discovery



- Méta-données (valeurs champ)
 - Signification business
 - Type
 - Longueur/précision
 - Valeurs valides (liste/range)
 - Null autorisé
 - Unique
 - Séquence
 - Formatage
 - Conventions d'encodage (null, etc)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

28



Data Profiling Discovery



- Analyse automatique
- Niveau de complexité variable
- Détermine **ce qui est**, pas ce qui est possible ou ce qui est permis.
- Méta-données (valeurs attribut)
 - Signification business
 - Type
 - Longueur/précision
 - Valeurs valides (liste/range)
 - Null autorisé
 - Unique
 - Séquence
 - Formatage
 - Conventions d'encodage (null, etc)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

29



Data Profiling Discovery



- Discovery (*Attribute properties*)
 - Fréquence
 - Patterns
 - Types de données effectifs (date, entiers)
 - Statistiques descriptives (min, max, médiane, moyenne, écart-type, distribution).
 - Nombre de *null*
 - Outliers
- Possibilité de "*drill down*"

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

30



Data Profiling Exemples



"Loan"

From <http://www.dataflux.com>

METRIC NAME	METRIC VALUE
Data Type	double
Primary Key Candidate	no
Unique Count	1140
Uniqueness	70.11
Pattern Count	(not applicable)
Minimum Value	-223000
Maximum Value	9999999
Minimum Length	(not applicable)
Maximum Length	(not applicable)
Null Count	2
Blank Count	(not applicable)
Actual Type	double
Count	1628
Data Length	53 bit
Mean	114348.170972
Median	4888499.5
Mode	0
Non-Null Count	1626
Nullable	YES
Ordinal Position	7
Decimal Places	0
Standard Deviation	429438.361236
Standard Error	10649.778281

"Phone number" Pattern analysis

PATTERN	COUNT	PERCENTAGE
999-999-9999	3166	96.73
(999)999-9999	42	1.28
(999) 999-9999	34	1.04
999 99 9999 999	20	0.61
999 999 9999	5	0.15
999-999-AAAA	2	0.06
9-999-999-9999	2	0.06
a	1	0.03
99 99 9999 999	1	0.03

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

31

Data Profile Manager: HR_PROFILE

Profile Edit View Window Help

Object Trees: HR_PROFILE, XWEEK, HRCOUNTRIES, HRDEPARTMENTS, **HREMPLOYEES**, HRJOBS, HRJOB_HISTORY, HRLOCATIONS, HRREGIONS

Property Inspector: Null Value Representation: null, Random Sample Rate: 100, Copy Data Into Workspace: [checked]

PROFILE CONFIGURATION

- Domain Value Compliance Min: 10
- Minimum Relationship Percent: 75
- Enable Relationship Discovery: [checked]
- Enable Pattern Discovery: [unchecked]
- Left Hand Side Attributes: 1
- Domain Value Compliance Min: 2
- Minimum Redundancy Percent: 100
- Domain Discovery Max Distinct: 50
- Enable Domain Discovery: [checked]
- Enable Redundant Columns Discovery: [checked]
- Enable Unique Key Discovery: [checked]
- Domain Discovery Max Distinct: 100

Profile Results Canvas: The following are the domain analysis results for HREMPLOYEES, which has 11 columns and 107 rows.

Columns	Found Domain	% Compliant	Stir-Sigma
HIRE_DATE		0%	10001.50
JOB_ID	ST_CLERK SA_REP SH_CLERK	65%	10001.50
LAST_NAME		0%	10001.50
MANAGER_ID	100	13%	10001.50
PHONE_NUMB...		0%	10001.50
SALARY		0%	10001.50

Data Drill Panel: Here are drill results on HREMPLOYEES column JOB_ID, related to Domains.

Distinct values: All

JOB_ID	# Rows	% of 107
8	1	1%
9	1	1%
10	1	1%
11	2	2%
12	5	5%
13	5	5%
14	5	5%
15	5	5%
16	5	5%
17	20	19%
18	20	19%
19	30	28%

Rows for the selected distinct value:

COMMISSION...	DEPARTMEN...	EMAIL	EMPLOYEE_ID	FIRST...
1	50	AWALSH	196	Alana
2	50	ABULL	185	Alexis
3	50	ACABRO	187	Arthor
4	50	BEVERETT	193	Britney
5	50	DOCONNEL	198	Donald
6	50	DGRANT	199	Douglas
7	50	GGEON	183	Girard
8	50	JFLEAUR	181	Jean
9	50	JOLLY	189	Jennife
10	50	JOELLING	186	Julia
11	50	KCHLING	188	Kelly

Displaying 19 Rows out of 19 | more

Displaying 20 Rows out of 20 | more

Oracle Business Intelligence logo



Data Profiling

Indices de problèmes



- Sur base de l'analyse automatique,
 - Attributs non utilisés ou peu utilisés
 - Représentations incohérentes
 - Représentations de NULL (vide, N/A, etc).
 - Valeurs inutilisées
 - Outliers (~ valeur inhabituellement grande ou petite)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

33



Data Profiling

Intervention humaine

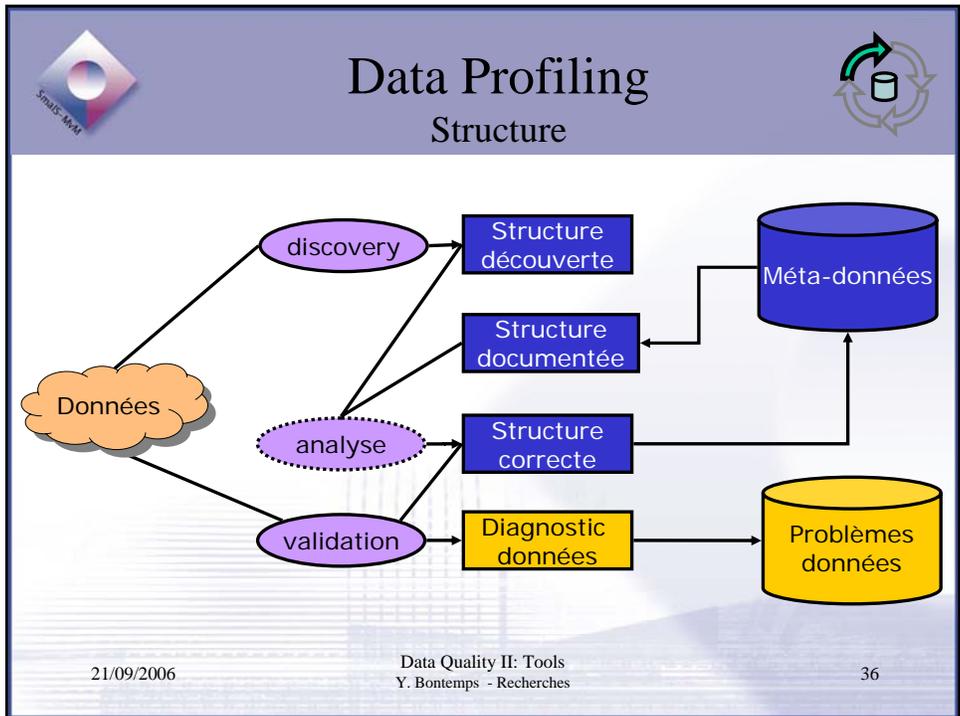
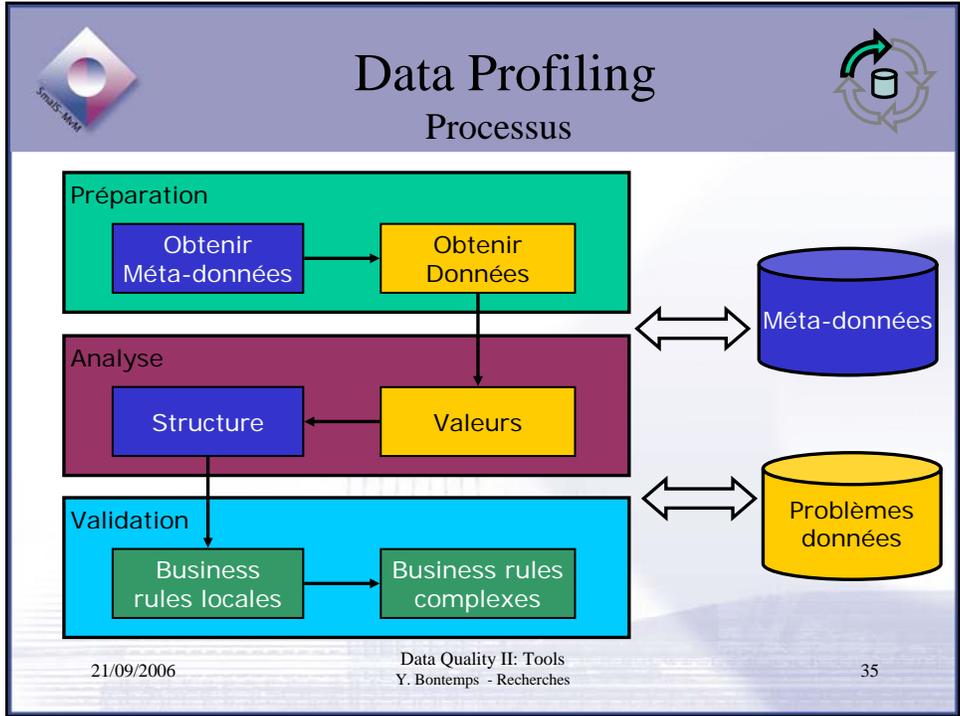


- Inspection visuelle
 - Valeurs extrêmes
 - Caractères spéciaux
 - Random walks
- Signification "business"
- Valeurs autorisées

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

34





Data Profiling Structurel



- Propriétés:
 - Identifiants
 - ONSS_NR
 - Relations (*Jointures*)
 - FORM_JUR est un code documenté dans la table annexe FORMES_JURIDIQUES.
 - Dépendances fonctionnelles
 - INS_CODE → VILLE
 - Synonymes

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

37



Data Profiling Besoin d'un outil ?



- Méthode ad-hoc (SQL, ...) ?
 - Effort nécessaire (écriture requêtes)
 - Manque de support méthodologique
 - Risque de manquer des informations importantes
 - Certaines analyses pas possibles:
 - Dépendances fonctionnelles
 - Synonymes
 - Drill-down
- Outils dédiés dépassent ces limitations

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

38



Data Profiling Reference Book



- **Data Quality: the Accuracy Dimension**, *Jack Olson*

Elsevier, 2002, The Morgan-Kaufmann Series in Database Management.



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

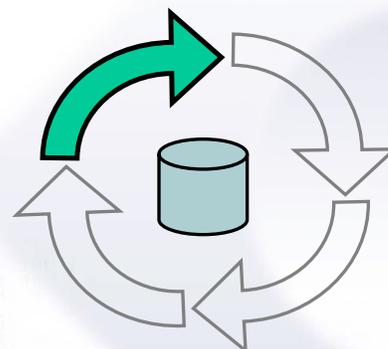
39



Data Profiling Fin du processus



- **Repositories:**
 - Méta-données
 - Colonnes
 - Structure
 - Règles
 - Contenu
 - Distribution
 - ...
 - Problèmes de données
 - Colonnes
 - Structure
 - Règles



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

40



Data Profiling

Outils commerciaux



- Points d'attention
 - Support méthodologique et collaboration
 - Ouverture du repository
 - Fonctionnalités de "discovery"
 - Expressivité du moteur de règles
(lang. contraintes built-in, reg-ex, SQL, javascript, etc)
 - Performance (!) et gestion des jobs
 - Intégration avec phases ultérieures de DQ
 - Possibilités de trends analysis (monitoring)

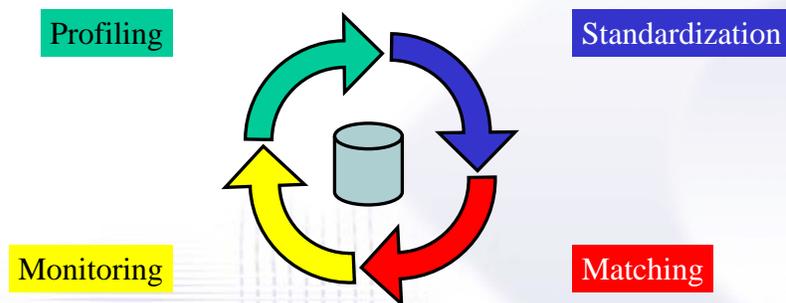
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

41



Cycle d'utilisation



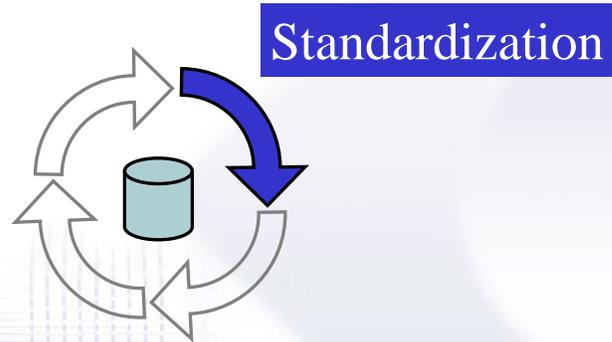
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

42



Cycle d'utilisation



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

43



Standardisation



- Buts
 - Cohérence des conventions de représentation
 - Correction de certains problèmes identifiés par le profiling
 - Enrichissement des données
 - Structuration de champs non-structurés
 - Utilisation de look-up tables

21/09/2006

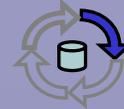
Data Quality II: Tools
Y. Bontemps - Recherches

44



Standardisation

Conventions de représentation



- Domaines restreints
 - Définition d'une transformation
Dom Invalide → Dom Valide
 - Directement dans outil de Profiling
- Domaines plus complexes (noms, adresses, etc)
 - Parsing + enrichissement
 - Domain-specific information
 - Données de références externes

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

45



Standardisation

Domaines Complexes



781114-269.56 | Yves Bontemps | Rue Prince Royal 102 Bruxelles

↓ Parsing

781114-269.56 | Yves | Bontemps | Rue Prince Royal | 102 | Bruxelles

↓ Enrichissement

781114-269.56 | Yves | Bontemps | Rue **du** Prince Royal | 102 | **Ixelles**

Male → 269 → Koninklijke Prinsstraat → Ixelles → Elsene → 1050

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

46



Standardisation

Domaines Complexes



- Programmation
 - Parsing

- Enrichissement

- "Encodage" de la connaissance du domaine

```

public class PersonStandardiser {
    protected List<String> decomposeNN(String nationalNumber) {
        List<String> decomposition = new ArrayList<String>(5);
        assert(nationalNumber.length() == 14);
        decomposition.add(nationalNumber.substring(0,2));
        decomposition.add(nationalNumber.substring(2,4));
        decomposition.add(nationalNumber.substring(4,6));
        decomposition.add(nationalNumber.substring(7,10));
        decomposition.add(nationalNumber.substring(11,13));
        return decomposition;
    }
    protected void enrich(Person pers, List<String> decomposedNN) {
        if (Integer.parseInt(decomposedNN.get(3)) % 2 == 1) {
            pers.setGender(Sex.MALE);
        }
        else {
            pers.setGender(Sex.FEMALE);
        }
        pers.setBirthDate(
            decomposedNN.get(2)
            + "/" + decomposedNN.get(1)
            + "/" + decomposedNN.get(0)
        );
    }
}

```



Standardisation

Parsing



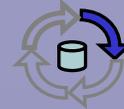
- Patterns, grammaires, expressions régulières

Pattern: 99.99.99-999.99
 Décomp: AN M J NR CD
 Règle enrich: NR mod 2 = 0 → sexe = 'M'
 NR mod 2 = 1 → sexe = 'F'
 DateNaiss = J '/' M '/' AN



Standardisation

Parsing – Outils commerciaux



- Base de connaissance
 - Corpus de règles
 - Gigantesque (50 000+ règles)
 - Années d'expérience
- Points d'attention:
 - Flexibilité/Exceptions
 - Fréquence des mises à jour
 - Adaptation des règles au contexte particulier (pays, région, langue, culture)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

49



Standardisation

Parsing - Futur



- Techniques d'Intelligence Artificielle
- Modèles probabilistes (Hidden Markov Models)
- Apprentissage → adaptation au contexte.
- Sur le marché dans 2 à 3 ans.

21/09/2006

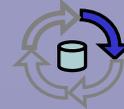
Data Quality II: Tools
Y. Bontemps - Recherches

50



Standardisation

Etapes



- Enrichissement
 - Calculs
 - Look-up tables
 - Enterprise Reference Data (Master Data)
 - Tables "annexes" (ex: codes postaux, etc).
 - Outil commercial (géographique, démographique, etc).
 - Information annexe
 - "Adresse invalide", "N° de maison inexistant" p. ex.

21/09/2006

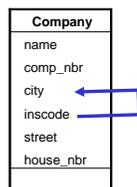
Data Quality II: Tools
Y. Bontemps - Recherches

51

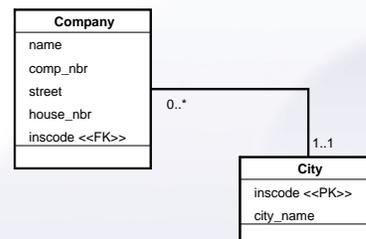


Standardisation

Structure



Dépendance
fonctionnelle



Normalisation

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

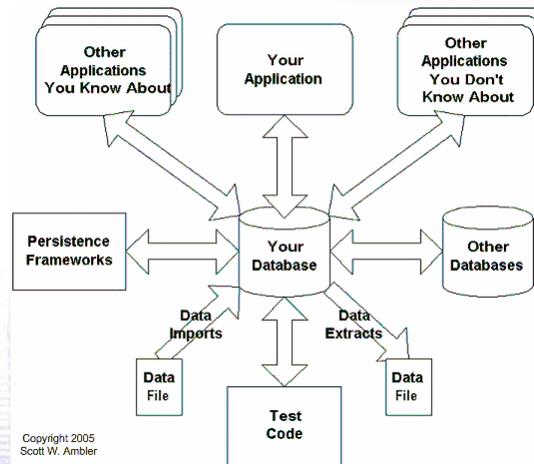
52



Standardisation Structure



- Modification du schéma DB opérationnelle
- Difficile à mettre en œuvre car couplage fort
→ co-évolution
- Pas de support logiciel spécifique



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

53



Standardisation Résumé



- Assurer la cohérence des représentations
- 2 objectifs:
 - Préparation du matching
 - Amélioration de la base de donnée source (en continu)
- Outils commerciaux laissent la structure intacte

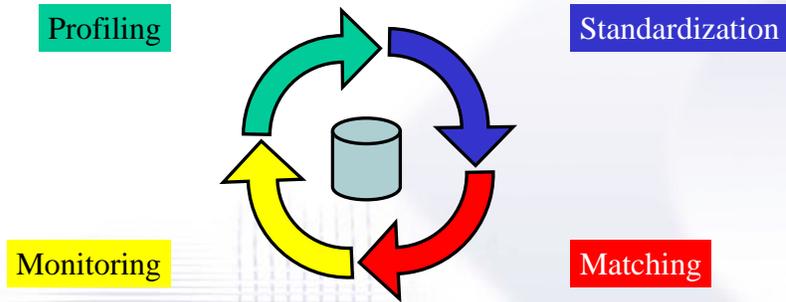
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

54



Cycle d'utilisation



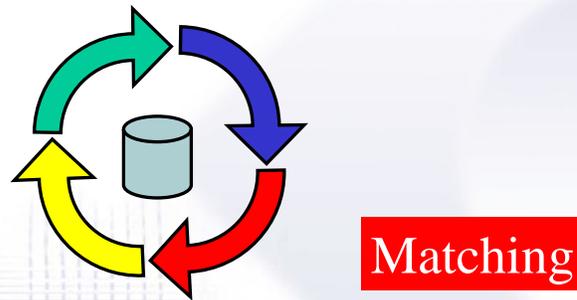
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

55



Cycle d'utilisation



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

56



Matching



- Déduplication
 - Détection des doublons
- Matching
 - Eviter doublons (ex: Data Integration)
 - Relier des BD → déduire de l'information, détecter des incohérences

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

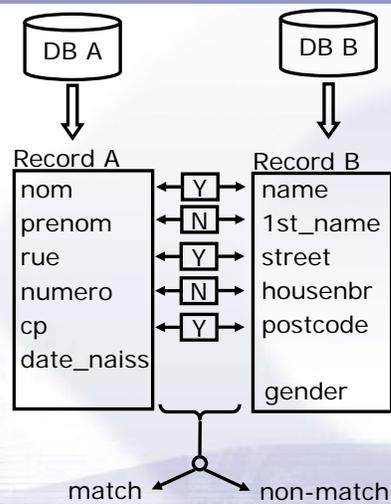
57



Matching



- Champ par champ
- Par record
- Remarques
 - Performance
 - Outils commerciaux



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

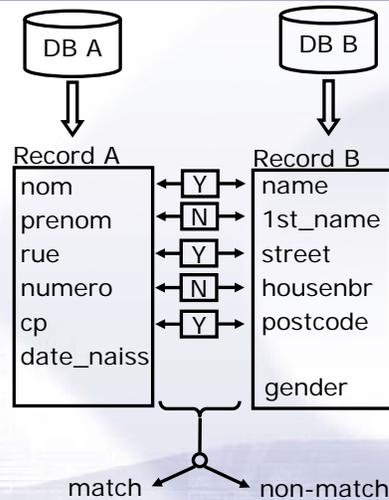
58



Matching



- **Champ par champ**
- Par record
- Remarques
 - Performance
 - Outils commerciaux



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

59



Matching

Champ par champ



- Chaînes de caractères
(noms, noms de rue, ...)
- Deux champs sont "les mêmes"
 - AFSCA "=" Agence Fédérale pour la Sécurité de la Chaîne Alimentaire
 - Yves "=" Yves
 - Yves "=" yves
 - Yves "=" Yvse
 - Bontemps "=" Bontan
 - Bontemps, Yves "=" Yves Bontemps

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

60



Matching

Familles de comparaisons



Booléennes	<p>Prédicats et règles</p> <ul style="list-style-type: none"> Egale, Suffixe de Stems from <p><i>Boulangier ~ Boulangerie</i> <i>S.A. ~ Société Anonyme</i></p>	<p>Phonétiques</p> <ul style="list-style-type: none"> Soundex Metaphone <p><i>Bontant ~ Bontemps</i></p>
	<p>Word-based</p> <ul style="list-style-type: none"> Edit distances Lettres communes <p><i>Bontemps ~ Bnotmps</i></p>	<p>Token-based</p> <ul style="list-style-type: none"> Cosinus Recursive <p><i>Office National des matières fissiles ~ Office des matières fissiles</i></p>
Similarité		

21/09/2006 Data Quality II: Tools 61
Y. Bontemps - Recherches



Matching

Familles de comparaisons



Booléennes	<p>Prédicats et règles</p> <ul style="list-style-type: none"> Egale, Suffixe de Stems from <p><i>Boulangier ~ Boulangerie</i> <i>S.A. ~ Société Anonyme</i></p>	<p>Phonétiques</p> <ul style="list-style-type: none"> Soundex Metaphone <p><i>Bontant ~ Bontemps</i></p>
	<p>Word-based</p> <ul style="list-style-type: none"> Edit distances Lettres communes <p><i>Bontemps ~ Bnotmps</i></p>	<p>Token-based</p> <ul style="list-style-type: none"> Cosinus Recursive <p><i>Office National des matières fissiles ~ Office des matières fissiles</i></p>
Similarité		

21/09/2006 Data Quality II: Tools 62
Y. Bontemps - Recherches



Data Matching

Comparaisons booléennes



- Relations booléennes
 - is-equal-to
 - is-a-prefix-of
 - is-a-suffix-of
 - is-an-abbreviation-of
 - stems-from
 - ...
- Exemples:
 - Bontemps ~ Bontemps
 - Bont ~ Bontemps
 - Emps ~ Bontemps
 - S.A. ~ Soc. Anon.
 - Bakkerij ~ Bakker

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

63



Matching

Rule-based



- Règles codées à la main
 - Hard-wired (Java, C, COBOL, ...)
 - Rule engine (declarative rules: prédéfinie et user-defined)
- Peuvent être fournies dans un outil commercial (black-box ?)
- ☺ Meilleure utilisation de la connaissance du domaine
- ☹ Fastidieux
- ☹ Coûteux (à créer/maintenir)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

64



Matching

Familles de comparaisons



Booléennes

Prédicats et règles

- Egale,
- Suffixe de
- Stems from

Boulangier ~ Boulangerie
S.A. ~ Société Anonyme

Phonétiques

- Soundex
- Metaphone

Bontant ~ Bontemps

Similarité

Word-based

- Edit distances
- Lettres communes

Bontemps ~ Bnotmps

Token-based

- Cosinus
- Recursives

*Office National des
matières fissiles ~ Office
des matières fissiles*

21/09/2006

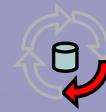
Data Quality II: Tools
Y. Bontemps - Recherches

65



Data Matching

Phonétique



- Erreurs de retranscription de l'oral vers l'écrit.
- Ex: Bontemps:
 - Montand
 - Bontant
 - Bonton
 - Bontand
 - Bautemps
 - ...

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

66



Data Matching Phonétique



- Principe:
 - Mot $m \rightarrow$ représentation de sa prononciation $p(m)$.
 - $p(m) = p(m') \Leftrightarrow m \sim m'$
- cfr indexation par hashing

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

67



Data Matching Phonétique



- Russel Soundex Algorithm (1918)
 - Garder première lettre
 - Supprimer a,e,h,i,o,u,w,y
 - Codage:
 - 1: B,F,P,V
 - 2: C,G,J,K,Q,S,X
 - 3: D,T
 - 4: L
 - 5: M, N
 - 6: R
 - Garder 4 premiers symboles (padding avec 0)
- Exemple
 - Bontemps
BNTMPS
B53512
B535
 - Bondant
BNDNT
B5353
B535

21/09/2006

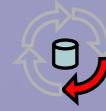
Data Quality II: Tools
Y. Bontemps - Recherches

68



Data Matching

Phonétique



- Variantes:
 - Metaphone et Double Metaphone
 - NYSIIS et NameSearch®
 - Daitch-Mokotoff
 - Langues Slaves et Germaniques
 - 54 entrées
 - Fonem
 - Français
 - 64 règles pour mettre les règles en forme normale
 - Phonex, etc...

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

69



Matching

Familles de comparaisons



Booléennes	Prédicats et règles <ul style="list-style-type: none">• Egale,• Suffixe de• Stems from <p><i>Boulangier ~ Boulangerie</i> <i>S.A. ~ Société Anonyme</i></p>	Phonétiques <ul style="list-style-type: none">• Soundex• Metaphone <p><i>Bontant ~ Bontemps</i></p>
	Similarité	Word-based <ul style="list-style-type: none">• Edit distances• Lettres communes <p><i>Bontemps ~ Bnotmps</i></p>

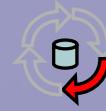
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

70



Data Matching Similarité



- Remplacer l'approche booléenne (0-1) par une approche plus fine
- Idée:
 - Similarité variable entre
 - Bontemps,
 - Smith,
 - Potnamp,
 - B0tenps.

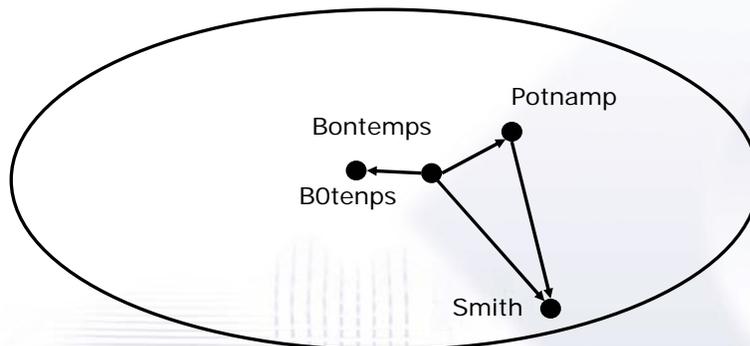
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

71



Data Matching Similarité



21/09/2006

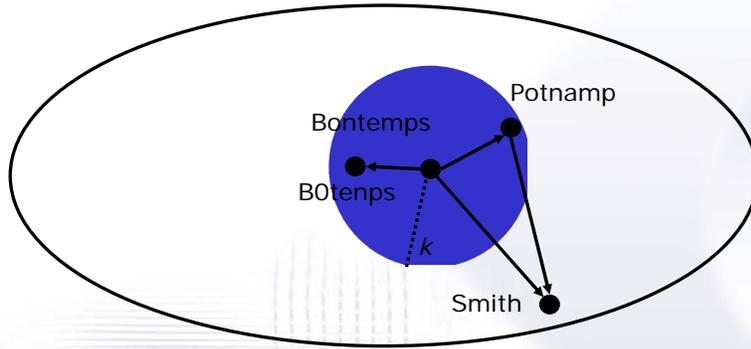
Data Quality II: Tools
Y. Bontemps - Recherches

72



Data Matching

Similarité



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

73



Matching

Familles de comparaisons



	<p>Prédicats et règles</p> <ul style="list-style-type: none"> Egale, Suffixe de Stems from <p><i>Boulangier ~ Boulangerie</i> <i>S.A. ~ Société Anonyme</i></p>	<p>Phonétiques</p> <ul style="list-style-type: none"> Soundex Metaphone <p><i>Bontant ~ Bontemps</i></p>
Booléennes		
	<p>Word-based</p> <ul style="list-style-type: none"> Edit distances Lettres communes <p><i>Bontemps ~ Bnotmps</i></p>	<p>Token-based</p> <ul style="list-style-type: none"> Cosinus Recursive <p><i>Office National des matières fissiles ~ Office des matières fissiles</i></p>
Similarité		

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

74



Matching

Edit distance



- Distance entre deux mots = nombre de "fautes" à effectuer pour passer de l'un à l'autre
 - Insertion (I)
 - Effacement (D)
 - Substitution (S)
- Exemple
 - Bontemps
 - Bontemps (I)
 - Bnontemps (D)
 - Bnontemps (D)

→ 3 opérations
- Coût par faute
 - OCR, Clavier, etc.

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

75



Data Matching

Lettres communes



- Jaro:
 - Fenêtre de $(m/2) - 1$
- Exemple
 - 3 caractères communs
 - 1 transposition
 - $d = \frac{1}{3} * \frac{3}{7} + \frac{1}{3} * \frac{3}{7} + \frac{1}{3} * \frac{(3 - 1)}{3}$
 - = 0.508

	S	L	U	I	T	E	N
S	1						
T					1		
R							
U			1				
D							
E						1	
L							

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

76



Matching

Familles de comparaisons



Booléennes	<div style="border: 1px solid black; background-color: #e6e6fa; padding: 5px; margin-bottom: 5px;"> <p style="text-align: center; margin: 0;">Prédicats et règles</p> <ul style="list-style-type: none"> Egale, Suffixe de Stems from <p style="margin: 0;"><i>Boulangier ~ Boulangerie</i> <i>S.A. ~ Société Anonyme</i></p> </div> <div style="border: 1px solid black; background-color: #e6e6fa; padding: 5px; margin-bottom: 5px;"> <p style="text-align: center; margin: 0;">Phonétiques</p> <ul style="list-style-type: none"> Soundex Metaphone <p style="margin: 0;"><i>Bontant ~ Bontemps</i></p> </div>	
Similarité	<div style="border: 1px solid black; background-color: #d8bfd8; padding: 5px; margin-bottom: 5px;"> <p style="text-align: center; margin: 0;">Word-based</p> <ul style="list-style-type: none"> Edit distances Lettres communes <p style="margin: 0;"><i>Bontemps ~ Bnotmps</i></p> </div> <div style="border: 1px solid black; background-color: #800040; color: white; padding: 5px;"> <p style="text-align: center; margin: 0;">Token-based</p> <ul style="list-style-type: none"> Cosinus Recursive <p style="margin: 0;"><i>Office National des matières fissiles ~ Office des matières fissiles</i></p> </div>	

21/09/2006
Data Quality II: Tools
Y. Bontemps - Recherches
77



Data Matching

Token-based



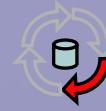
- Plusieurs éléments dans un record
- Exemples
 - Yves Bontemps
 - Organisme national des déchets radioactifs et des matières fissiles.

21/09/2006
Data Quality II: Tools
Y. Bontemps - Recherches
78



Data Matching

Token-based metrics



- Ignorer l'ordre des mots
- Métrique de Jaccard
$$\frac{|S \cap T|}{|S \cup T|}$$
- Bontemps, Yves
=?=
Yves Bontemps
- { { Bontemps; Yves } }
{ { Yves; Bontemps } }
- $2/2 = 1$

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

79



Data Matching

Token-based



- Term Frequency / Inverse Document Frequency (TF-IDF)
- Rareté d'un mot → importance.
- Exemples:
 - Organisme National des Déchets Radioactifs et des Matières Fissiles
 - Office des Matières Fissiles et Déchets Radioactifs

21/09/2006

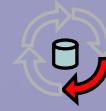
Data Quality II: Tools
Y. Bontemps - Recherches

80



Data Matching

Token-Based



- Approches récursives:
 - Utiliser une technique "word-based" (Levenshtein, etc) pour déterminer si 2 tokens sont les mêmes (*Office ~ Office*)
 - Utiliser une technique "token-based" pour combiner ces résultats
 - Réf: Soft TF-IDF, etc.

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

81



Matching

Familles de comparaisons



Booléennes	Prédicats et règles <ul style="list-style-type: none">• Egale,• Suffixe de• Stems from <i>Boulangier ~ Boulangerie</i> <i>S.A. ~ Société Anonyme</i>	Phonétiques <ul style="list-style-type: none">• Soundex• Metaphone <i>Bontant ~ Bontemps</i>
	Similarité	Word-based <ul style="list-style-type: none">• Edit distances• Lettres communes <i>Bontemps ~ Bnotmps</i>

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

82



Data Matching

En pratique



- Comparaison phonétique (KBO):
 - Décomposition en mots
 - Mots mis sous forme phonétique.
 - Variante de Soundex
 - Ignore termes habituels (S.A., ...)
 - Stemming
 - Comparaison (prise en compte de l'ordre des mots)

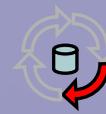
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

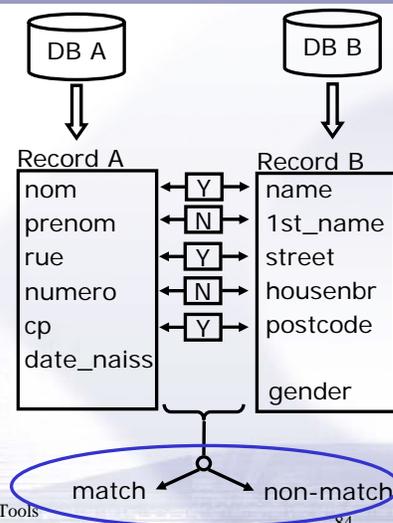
83



Matching



- Champ par champ
- **Par record**
- Performance



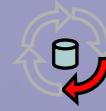
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

84



Data Matching Par record



- Approches:
 - Déterministe
 - Probabiliste (*global scoring*)

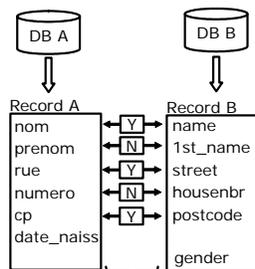
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

85



Data Matching Approche déterministe



	Nom	Prénom	Rue	Num	Cp	Result
1	Y	Y	-	-	Y	Match
2	N	-	-	-	-	No
3	Y	N	Y	Y	Y	Match
4	Y	N	Y	N	Y	Maybe
4	-	-	-	-	-	No

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

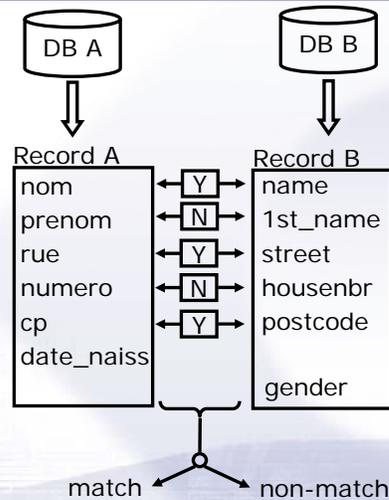
86



Matching



- Champ par champ
- Par record
- **Performance**



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

89



Data Matching Performances



- Matching = comparer tous les records de A avec tous les records de B
- Exemple:
 - A : 100 000 rec.
 - B : 100 000 rec.
 - ➔ 10 000 000 000 de comparaisons.
- ➔ Optimisations nécessaires...

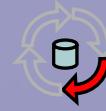
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

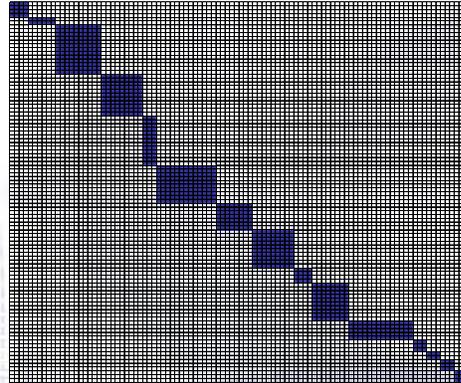
90



Data Matching Performances



- Blocking:
 - Donnée critique → former des blocs
 - Les records sont comparés bloc par bloc.
 - Ex: Code Postal, ...
 - Trade-off précision et performance
 - Nécessite indexation/tri.
 - Multi-passes possible
- Sliding window, etc



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

91



Data Matching Outils commerciaux



- Comparaison champ-par-champ
 - Best-of-breed, best-of-market (black-box ?)
- Agrégation
 - Global scoring (dans la plupart des cas)
 - Déterministe
- Performances
 - Blocking inclus (obligatoire)
 - DBMS extérieur ou gestion propre.
- Mode online et mode batch.

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

92



Cycle d'utilisation

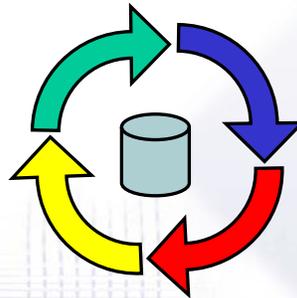


Profiling

Standardization

Monitoring

Matching



21/09/2006

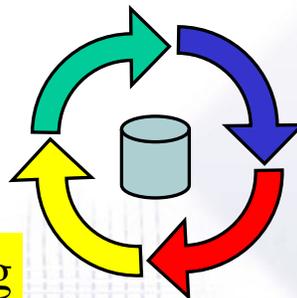
Data Quality II: Tools
Y. Bontemps - Recherches

93



Cycle d'utilisation

Monitoring



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

94



Monitoring



- Surveillance de l'évolution de la qualité
 - Régulière
 - Planifiée
 - Automatisée
- Tableaux de bord et alarmes
- Trends analysis

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

95



Monitoring



- Profiling tools
 - Violation des règles business (méta-données prescriptives)
 - Evolution des fréquences de valeur
 - Ex: +150% de N/A, etc.
- Possibilité de stocker des profils d'analyse et de les répéter au cours du temps

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

96



Monitoring



- Aspect nécessaire,
- Data Quality = Processus continu

- Pertinence des outils
 - Intégration
 - Architecture IT (impact sur performances).
 - Decision Support Systems
 - Fonctionnalités proposées (faibles)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

97



- Introduction
- Outils de Data Quality: concepts
 - Data Profiling
 - Standardisation
 - Matching
 - Monitoring
- **Outils commerciaux**
 - Case study
 - Architecture
- Conclusion

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

98



Data Quality Marché

- Fournisseurs principaux
- Case studies/Proof of concept
- Offre logicielle

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

99



Fournisseurs principaux

- Avril 2006: Magic Quadrant Gartner → marché reconnu.
- Deux phases de consolidation
 - Profiling → DQ
 - DQ → BI/ETL
- Futur: DQ = commodity?



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

100



Outils commerciaux Offre

- Data Profiling
 - Discovery
 - Assertive testing
 - Collaboration (via email)
 - Langage de règles variable
- (Monitoring)
- Standardisation et matching
 - Base de connaissance (règles, adresses, ...)
 - Performance
 - Templates de projet
 - Connecteurs
 - Updates



21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

101



Case studies

- Répertoire des employeurs ONSS
- Détection de doublons
 - Insertion des doubles immatriculations
- 230 000 entrées et 10 colonnes analysées, dont
 - Dénomination,
 - Adresse
 - Forme juridique, ...

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

102



Case studies

- "Workshops" d'une journée
- Planning
 - Profiling (2h)
 - Matching (4h)
 - Debriefing (1h30)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

103



Case studies Enseignements

- Performances
 - Rapidité de mise en œuvre (situation simple)
 - Performances (6 min. à 30 min. machine low-end)
 - Facilité d'utilisation (env. graphique, pas de programmation)
- Standardisation
 - Adresses bilingues: support variable.
 - Dénomination: pers. physiques vs pers. morales
- Matching
 - Champ par champ: black-box (pas un obstacle)
 - Record par record: Déterministe → traçabilité
 - Qualité des résultats

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

104



Outils commerciaux

- Ordres de prix (licences)
 - Profiling:
 - 25 000 € (1 named user)
 - 10 000 € par named user suppl.
 - Peu de users dans les entreprises
 - Standardisation, Matching:
 - 100 000 €

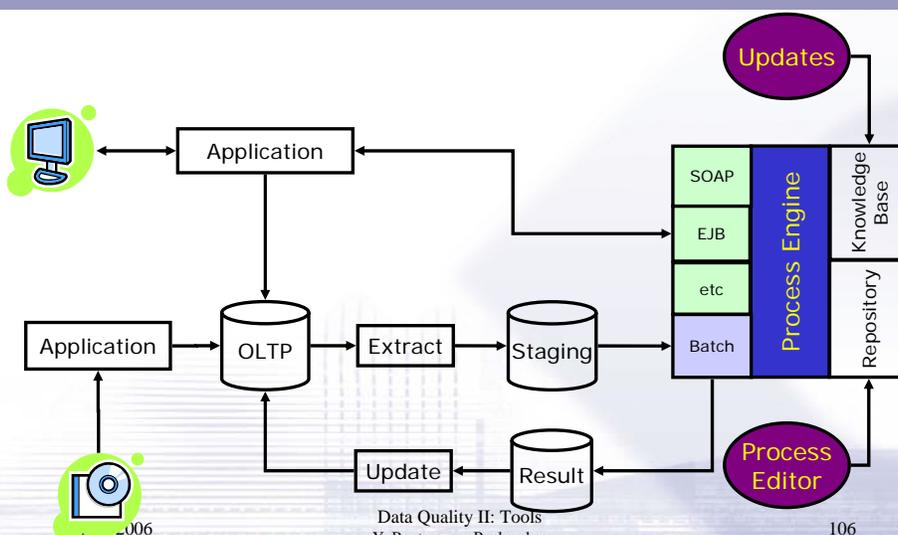
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

105



Architecture



Data Quality II: Tools
Y. Bontemps - Recherches

106



- Introduction
- Outils de Data Quality: concepts
 - Data Profiling
 - Standardisation
 - Matching
 - Monitoring
- Outils commerciaux
 - Case study
 - Architecture
- **Conclusion**

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

107



Conclusion

- Data Quality Tools ?
 - Fonctionnalités ?
 - Profiling
 - Standardisation/Enrichissement
 - Matching/Deduplication
 - **Monitoring**
 - Apport ?
 - Economie de l'expérience !
 - Performances

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

108



Conclusion

Scénarios

- Profiling
 - Au début d'un projet, extraction, documentation et vérification des données
 - Quelques jours (analyse)
 - Validation de toutes les hypothèses !

 - Diminution risques projet

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

109



Conclusion

Scénarios d'utilisation

- Standardisation/Déduplication batch
 - Création d'entités via Web ou batch

 - Déduplication en batch
 - Répétitif (mensuel, hebdomadaire, ...)
 - Traitement:
 - Surviving record
 - Gestion manuelle
 - Linkage des records (trace/info supplémentaire)

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

110



Conclusion

Scénarios d'utilisation

- Validation on-line
 - Lors de l'encodage,
 - Validation
 - Standardisation
 - Déduplication
 - Feed-back vers l'utilisateur
- Impact sur l'application "front-office"

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

111



Conclusion

Facettes d'une solution

Management

- Exec. DQ group
- Data Mgmt
- Data stewards
- ...

Information

- Meta-data
- Data Models
- Requirements
- Metrics, ...

Data as an Asset

Activités

- Root-Cause Analysis
- Data Integration
- Software Eng.
- Business Proc. Eng.
- Cleansing, etc.

Bénéfices

- Utilisateurs
- Développeurs
- Clients
- ...

21/09/2006

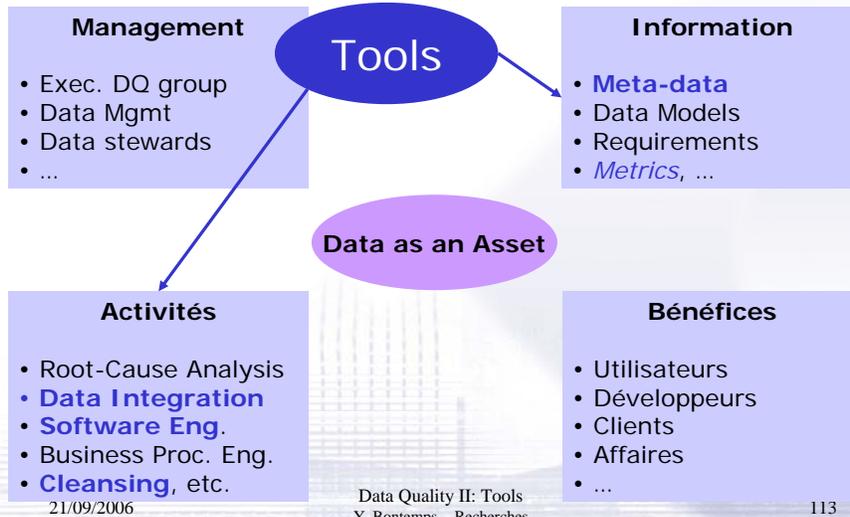
Data Quality II: Tools
Y. Bontemps - Recherches

112

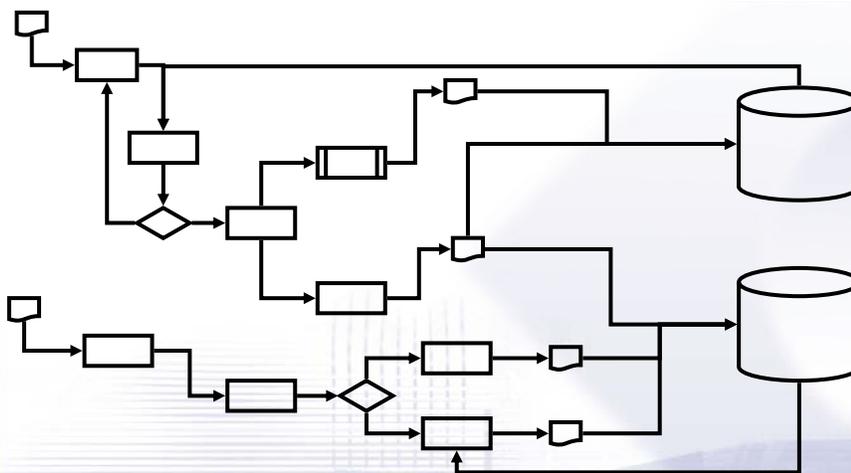


Conclusion

Facettes d'une solution



Data Quality



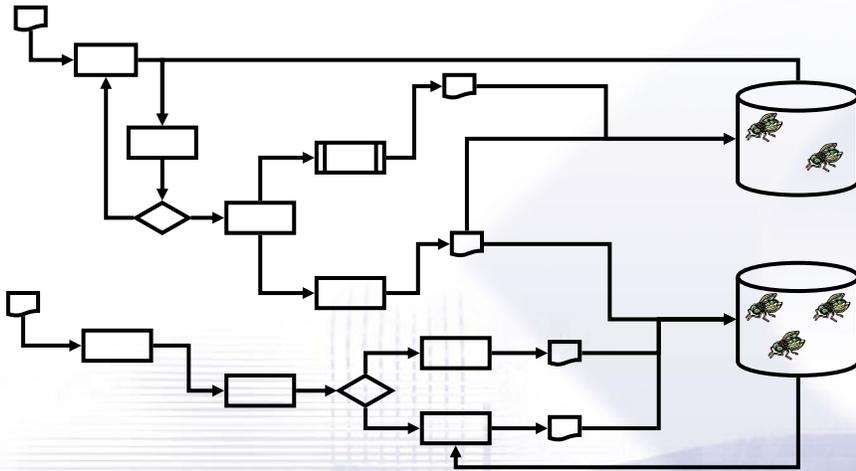
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

114



Data Quality Profiling



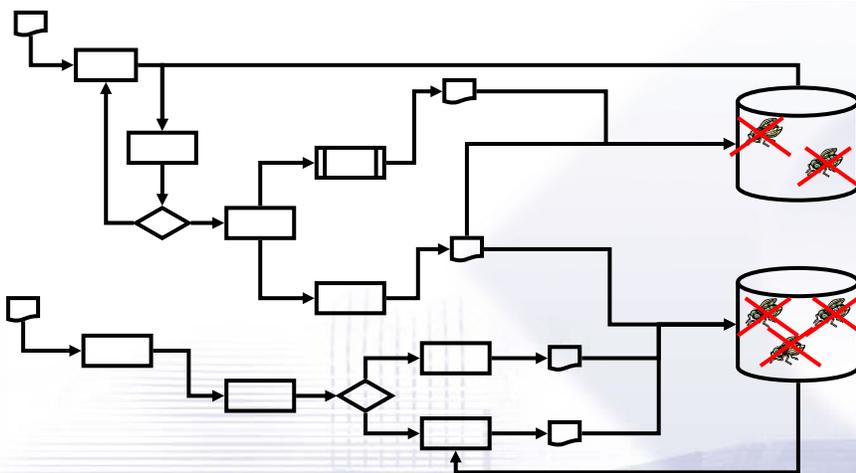
21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

115



Data Quality Standardisation & Matching

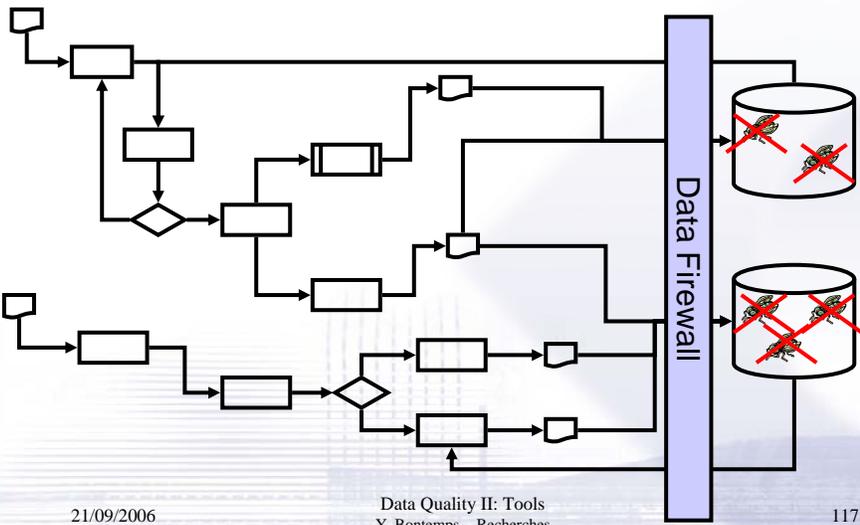


21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

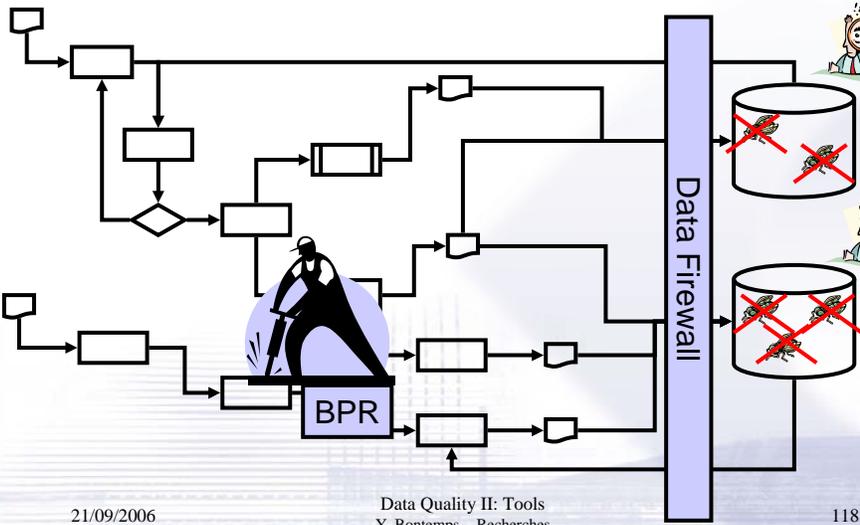
116

 **Data Quality Infrastructure**



21/09/2006 Data Quality II: Tools
Y. Bontemps - Recherches 117

 **Data Quality**



21/09/2006 Data Quality II: Tools
Y. Bontemps - Recherches 118



"Data Quality @ SmalS-MvM"

Cellule "data quality" (section "recherches")

- En collaboration avec les autres équipes de la société:
 - Sensibilisation à la qualité des données,
 - Formations,
 - Mise en place d'indicateurs,
 - Mise en place de groupes de travail & de suivi,
 - Actions spécifiques (root-cause analysis, etc),
 - Analyses de l'existant (impact, ...),
 - Aide à la mise en place d'outils,
- Etudes et publications de travaux
- Consultances au sein de l'administration fédérale belge

21/09/2006

Data Quality II: Tools
Y. Bontemps - Recherches

119