

Ceuppens, Cobben, Cobbin, Cober, Cobet, Cobie, Cobus,
 Cobussen, Coobs, Coops, Cop, Cöp, Copijn, Copini, Copius,
 Coppée, Coppen, Coppens, Coppes, Coppy, Coppin, Coppis,
 Coppus, Cops, Cup, Cuppé, Cuppen, Cuppens, Cuppes, Jaacke,
 Jacob, Jacobi, Jaakob, Jacobus, Jacobusse,
 Jacquemijns, Jakob, Jacobus, Jacobus, Jacobus, Jacobus,
 Kobben, Köbber, Köbber, Köbber, Köbber, Köbber,
 Kobs, Kobus, Kobus, Kobus, Kobus, Kobus, Kobus,
 Koops, Koopsen, Koopsen, Koopsen, Koopsen, Koopsen,
 Koopsma, Kop, Köpcke, Kopee, Kopjes, Köpcke, Kopp, Köpp,
 Koppe, Koppei, Köppel, Koppen, Köppen, Koppennens, Koppens,
 Koppers, Köppers, Koppes, Koppeij, Koppies, Koppj, Koppijn,
 Koppius, Koppj, Koppj, Koppj, Koppj, Koppj, Koppj,
 Kubin, Kubi, Kubus, Kubus, Kubus, Kubus, Kubus,
 Kuup, Jacob, Jacob, Jacob, Jacob, Jacob, Jacob,
 Yacob, Yacobi, Yacobus, Yacooob, Yacoub, Yacouba, Yacoubi,
 Yacoubian, Yacoubou, Yacubi, Yacubu, Yako, Yakob, Yakobchuk

NAMES

Gerrit Bloothoofd & David Onland, UIL-OTS Utrecht
 Martin Reynaert, TICC / Tilburg University
 Katrien Depuydt & Tanneke Schoonheim, INT Leiden

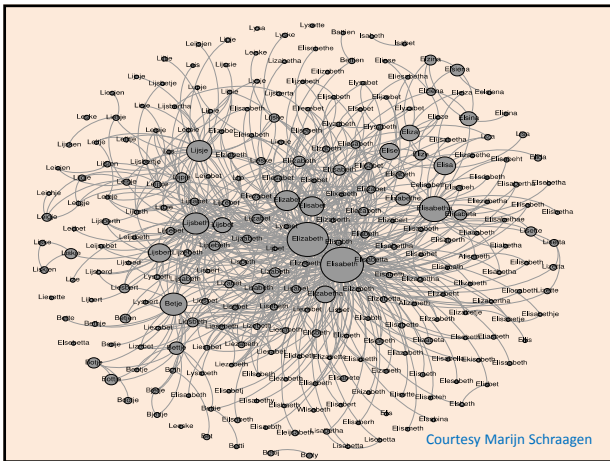
aims

standardization of 19th century person names for

- richer search results
- OCR post-processing
- nominal record linkage
- onomastic research

resulting NAMES corpus

in RDF format for Linked Open Data
& lexicon service



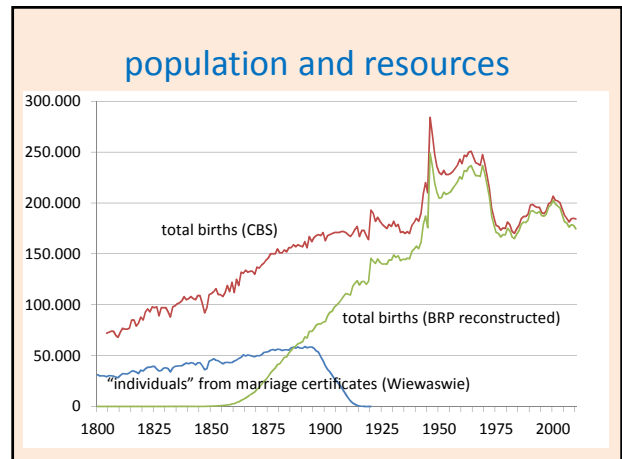
person name standardization

main issues

- on what basis?
 - onomastic expert knowledge
 - automatic learning from data
- data structure
 - handling ambiguities
- usage
 - search tools
 - nominal record linkage

automatic learning from data

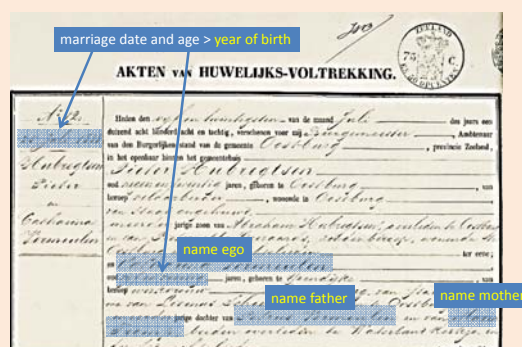
proven name variants
 from *identified* individuals in various documents
 followed by expert review



base material

564.000 surnames and 189.000 first names
 from 19th century sources
 (birth, marriage, death certificates)
 with 63 million person tokens
 from Catch LINKS project – Wiewaswie (2009 version)

certificate 3rd marriage



uniquely identifying information

- name ego, date and place of birth
- name ego, year of birth, names of parents
- name ego, year of birth, name of partner
-

factors determining uniqueness

- size of population (geographic spread)
- data completeness
- data accuracy
- name conventions
- social/family relationships

sufficient information

- first name ego
- three out of four names of parents
- one name may show up as variant (first name or surname) and variants can be harvested

proven-variant clustering

- surnames 127.154 into 15.114 standards
 compare:
 - corpus of Dutch surnames (CBG) 320.000 names, 59.000 in NAMES, 11.800 standards
 - surnames in Belgium and northern France (Debrabandere), 125.000 names, 49.000 in NAMES, 19.000 standards
- first names 42.196 into 2.900 standards (gender-dependent)
 compare:
 - first name dictionary (van der Schaar) 20.000 names, 12.500 in NAMES, 2.400 standards

Christoffel

CCHRISTOFFEL	CHRISTOFF	CHRISTOPHLE	CHRISTOFFEL
CGRISTOFFEL	CHRISTOFFE	CHRISTOPHOLUS	CHTSTOFFEL
CHRISTOFFEL	CHRISTOFFEL	CHRISTOPHONES	CRISTOFFEL
CHRISTOFFEL	CHRISTOFFELINA	CHRISTOPHONIS	CRISTOFFER
CHRISTOFF	CHRISTOFFELL	CHRISTOPHONUS	CRISTOP
CHRISTOPH	CHRISTOFFELLINA	CHRISTOPHORA	CRISTOFFEL
CHRISTOFF	CHRISTOFFER	CHRISTOPHORE	CRISTOFFER
CHRISTOFF	CHRISTOFFES	CHRISTOPHORES	CRISTOFFERUS
CHRISTOPHORUS	CHRISTOFFE	CHRISTOPHORI	CRISTOPH
CHRISTOPHORUS	CHRISTOPHE	CHRISTOPHORIS	CRISTOPHE
CHRISTOFFEL	CHRISTOFINA	CHRISTOPHORUS	CRISTOPHORUS
CHRISTOFFEL	CHRISTOFLE	CHRISTOPHOSUS	KIRSTOFFEL
CHRISTOFFEL	CHRISTOPORA	CHRISTOPHOZA	KIRSTOFFEL
CHRISTOFFELINA	CHRISTOPORIS	CHRISTOPHOZUS	KRISTOFFEL
CHRISTOFFER	CHRISTOPORUS	CHRISTOPHRE	KRISTOFFER
CHRISTOPH	CHRISTOPUS	CHRISTOPHORUS	STEFFERTIEN
CHRISTOPHORUS	CHRISTOPORA	CHRISTOPHORUS	STOFFEL
CHRISTEL	CHRISTOKE	CHRISTOPHUS	STOFFEL
CHRISTEPH	CHRISTOFFEL	CHRISTOPORUS	STVELIENA
CHRISTHOFFEL	CHRISTOP	CHRISTOTINA	STOFER
CHRISTHOP	CHRISTOFF	CHRISTPH	STOFFER
CHRISTOPH	CHRISTOFFEL	CHRISTPHEL	STOFFEL
CHRISTOPHORUS	CHRISTOPH	CHRISTPHORIS	STOFFEN
CHRISTOFFEL	CHRISTOPHARUS	CHRISTPHORUS	STOFFER
CHRISTOPHORUS	CHRISTOPHE	CHRISTOPHORUS	STOFFERT
CHRISTOPH	CHRISTOPHEL	CHRISTOPHOFEL	STOFFE
CHRISTOPH	CHRISTOPHELT	CHRISTOFFER	STOFFER
CHRISTOFFEL	CHRISTOPHER	CHRISTOFFEL	STOFFERS
CHRISTOFFELINA	CHRISTOPHERUS	CHRISTOPH	STOKKER
CHRISTOFFELINA	CHRISTOPHENA	CHRISTOFFEL	STOFFER
CHRISTOFFER	CHRISTOPHINA	CHROTOPEL	ZOFFER

Christiaans

CHRISTIAANSE	CHRISTIAANSE	CHTISTIAENS
CHRISTIAANS	CHRISTIAANSEN	CRISTIAAN
CHRISTIAAENS	CHRISTIAANSENS	CRISTIAANS
CHRISTIAEANS	CHRISTIAANSENS	CRISTIAANSE
CHRISTIAEANS	CHRISTIAANSENZEN	CRISTIAANSEN
CHRISTIAEANS	CHRISTIAEN	CRISTIAENS
CHRISTIAEANS	CHRISTIAENS	CRISTIAANS
CHRISTIAANS	CHRISTIAENSEN	KRISTIAANS
CHRISTIAANS	CHRISTIAENSENS	KRISTIAAN
CHRISTAANSE	CHRISTIAN	KRISTIAANS
CHRISTAEN	CHRISTIANA	KRISTIAANUS
CHRISTAENS	CHRISTIANE	KRISTIAANS
CHRISTAENSENS	CHRISTIANUS	KRISTIONS
CHRISTIANI	CHRISTIANUS	KRISTJAAN
CHRISTEAANS	CHRISTIAANS	KRISTJAANS
CHRISTIAAENS	CHRISTIAANSEN	
CHRISTIAAN	CHRISTIANUS	
CHRISTIAANS	CHRISTIAENS	

Luijk

L?CKEN	LUCKES	LUIKEN	LUKS
L?CKERS	LUCKS	LUIKENS	LUKTE
L?KEN	LUIBEN	LUIKER	LUKTE
LAEKENS	LUICKX	LUIKES	LUKTKEN
LAEKENS	LUIHEN	LUIKS	LUJKS
LAKENS	LUIBEN	LUIKT	LUX
LAKIUS	LUICK	LUIN	LUYCK
LALENS	LUICKEN	LUIX	LUYCKX
LOKES	LUICKU	LUK	LUYCK
LOUKENS	LUICKS	LUKA	LUYK
LUCA	LUICKSE	LUKE	LUYKE
LUCCA	LUICKX	LUKEN	LUYKEN
LUCK	LUICKX	LUKENS	LUYKS
LUCKA	LUICKX	LUKER	LUYKX
LUCKE	LUIK	LUKES	LUYN
LUCKEN	LUJK	LUKHEN	LÜCKE
LUCKENER	LUJKE	LUKIEN	LÜCKEN
LUCKER	LUIKEN	LUKIE	LÜCKER
LUCKERS	LUIKER	LUKKA	LÜCKERS
LUCKERT	LUIKES	LUKKEN	LÜKE
LUCKES	LUIKS	LUKKER	LÜKEN
LUCKNER	LUICKX	LUKKERS	LÜKS
LUCKS	LUJUM	LUKKERT	
LUCKX	LUJUN	LUKES	
LUDKEN	LUJK	LUKKIEN	
LUKEEN	LUJKE	LUKNER	

expert review

- surnames 15.114 standards into 11.539 standards
- first names 2.900 standards into 926 standards (gender-*in*dependent)

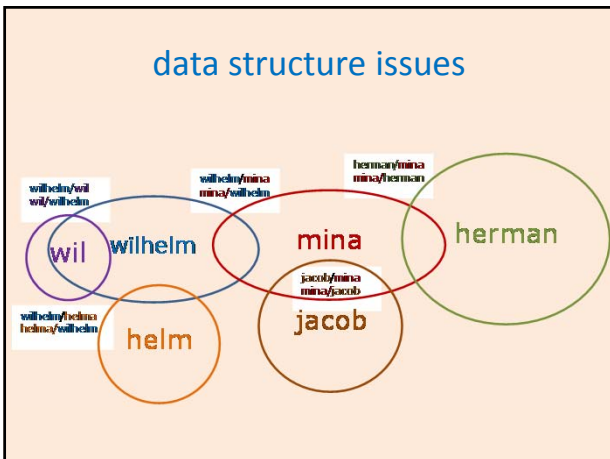
926 standards into ~500 base forms (first syllable etymology)

standards and application

choice of standards depends on application

- wide search: broad standards
- narrow search: fine standards

=> data structure for standards at different levels



full analysis

- learning variant properties (*TICLL, Text-Induced Corpus Clean up*)
 - context-dependent edit distance
 - clustered edit operations
- standardization of remaining names
 - 127.000 -> 564.000 surnames
 - 46.000 -> 189.000 first names

nominal record linkage



nominal record linkage

strategies:

pairwise edit-distance

- certificate-based or ego-based
- various (edit-)distance measures
- problems for large edit-distances
- quadratic computation problem, slow

sorted neighborhood

- ego-based
- **requires standardized data**
- matching within window (size=2) only
- fast

optimal field order for sorting key

for birth, marriage, death of ego

1. surname ego [distinguishes more than first name]
2. first name ego
3. sex ego
4. surname mother
5. first name mother
6. surname father [father not always known]
7. first name father
8. birth date ego [imprecise]

place of birth ego is not used [imprecise]

sorting example

ego surname	ego first name	sex	mother surname	mother first name	father surname	father first name	date birth (days)	role
doorn	wilhelmina	f	dishoek	petronella	doorn	abraham	657040	deceased
doorn	wilhelmina	f	dishoek	petronella	doorn	abraham	657103	bride
doorn	wilhelmina	f	doorn	jacoba	doorn	johannes	671562	child
doorn	wilhelmina	f	doorn	jacoba	doorn	johannes	671566	deceased
doorn	wilhelmina	f	nauta	geertrui	doorn	kornelis	647143	deceased
doorn	wilhelmina	f	schorer	adriana	doorn	hendrik	661241	child
doorn	wilhelmina	f	vlag	maria	doorn	leendert	697569	bride
doorn	wilhelmina	f	vlag	maria	doorn	leendert	697661	child

conclusion

- standardization based on proven variants
 - expert review needed
 - many uncertainties (short names)
 - flexible data structure required
 - extension to patronyms
-
- extremely useful for many applications
 - likely useful for older documents