# Processing and linking address data

Madis Nõmme Karl Kilgi

## Task description

#### Goal:

- build an address matching engine
- take addresses in textual form
- match them against structured DB
- compare results with official search

Solution available from github.com/madis/ses-2013

#### Data used

Real data from Maanteeamet (Estonian Road Administration) was used

- 2 datasets in CSV format:
  - 1. Tallinn 395,638 rows
  - 2. Whole Estonia 2,611,625 rows
- 37 fields for every row

	A1   T   W   J^A   ADOB_ID										
_4	A	В	C	D	E	F	G	Н	I	J	K
1	ADOB_ID	ADS_OID	ADOB_LIIK	ORIG_TUNNUS	ETAK_ID	ADS_KEHTIV	OLEK	ADR_ID	KOODAADRESS	TAISAADRESS	LAHIAADRESS AA
2	118006	LP00012285	LP	2119541		19.02.2004 06:40:40	K	108838	377840524000004DS00000000000000000	Harju maakond, Tallinna linn, NVµmme linnaosa, Kalevala tn	Kalevala tn K
3	121233	LP00016245	LP	2614698		28.05.2004 09:55:18	K	110991	377840596000004JP00000000000000000	Harju maakond, Tallinna linn, Pirita linnaosa, Sarapiku tee	Sarapiku tee K
4	121263	LP00016282	LP	2706927		28.05.2004 09:55:04	K	117726	377840176000004H300000000000000000	Harju maakond, Tallinna linn, Haabersti linnaosa, Kudu tn	Kudu tn K
5	121411	LP00016463	LP	2104757		19.02.2004 01:25:50	K	113229	377840387000004BU00000000000000000	Harju maakond, Tallinna linn, Lasnam√§e linnaosa, Tuulem√§e tn	Tuulemv§e tn K
6	121465	LP00016529	LP	2120362		19.02.2004 06:58:51	K	113777	377840524000003RC00000000000000000	Harju maakond, Tallinna linn, NVµmme linnaosa, Uku tn	Uku tn K
7	120020	LP00014757	LP	2105767		19.02.2004 01:52:14	K	115634	377840387000004BG00000000000000000	Harju maakond, Tallinna linn, Lasnamv§e linnaosa, VV§ike-Paala tn	VV§ike-Paala tn K
8	120241	LP00015028	LP	2114760		19.02.2004 04:42:09	K	108500	377840298000004AH00000000000000000	Harju maakond, Tallinna linn, Kesklinna linnaosa, Rahukohtu tn	Rahukohtu tn K
9	120247	LP00015035	LP	2120075		19.02.2004 06:52:54	K	116726	377840524000004E200000000000000000	Harju maakond, Tallinna linn, NVµmme linnaosa, Kasteheina tn	Kasteheina tn K
10	120270	LP00015063	LP	2119750		19.02.2004 06:46:20	K	108877	377840524000003Z900000000000000000	Harju maakond, Tallinna linn, NVµmme linnaosa, Puraviku tn	Puraviku tn K
11	118504	LP00012896	LP	2612376		28.05.2004 09:50:50	K	115414	377840298000005ZG0000000000000000	Harju maakond, Tallinna linn, Kesklinna linnaosa, Loite tn	Loite tn K
12	120568	LP00015428	LP	2114000		19.02.2004 04:22:56	K	112801	377840298000003KK00000000000000000	Harju maakond, Tallinna linn, Kesklinna linnaosa, Vana-Veerenni tn	Vana-Veerenni tn K
13	118732	LP00013177	LP	2107313		19.02.2004 02:21:08	K	108324	377840596000003S60000000000000000	Harju maakond, Tallinna linn, Pirita linnaosa, Kedriku tn	Kedriku tn K

## Importing & structuring data

With optimizations ~2000 rows/s Wrote a migration script

- took ~5 min to import ~400k rows
- data in relational DB

# Searching approaches

Simple query

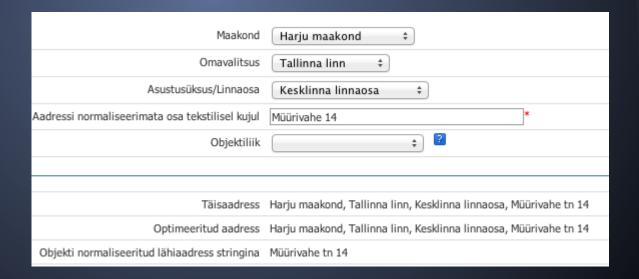
• SQL WHERE x LIKE '%...%'

Full text search -

- Sphinx
- Elasticsearch

## Validating results

Validated against http://xgis.maaamet.ee/adsavalik/ads Example queries



# Validating results

### Results

Jkn	ADR-	Objekti versiooni id	Liik	Tase 1	Tase 2	Tase 3	Tase 4	Tase 5	Tase 6	Tase 7	Tase	Koodaadress	normaliseerimata osa tekstilisel kujul
1	390906	3551857	mitteelukondlik hoone	Kood: 37 Nimi: Harju maakond Nimi liigisõnaga: Harju maakond	Kood: 784 Nimi: Tallinna linn Nimi liigisõnaga: Tallinna linn	Kood: 0298 Nimi: Kesklinna linnaosa Nimi liigisõnaga: Kesklinna linnaosa		Kood: 03QJ Nimi: Müürivahe tn Nimi liigisõnaga: Müürivahe tänav		Kood: 0PEN Nimi: 14 Nimi liigisõnaga: 14		377840298000003QJ00000PEN00000000	
2	390906	3061437	katastriüksus	Kood: 37 Nimi: Harju maakond Nimi liigisõnaga: Harju maakond	Kood: 784 Nimi: Tallinna linn Nimi liigisõnaga: Tallinna linn	Kood: 0298 Nimi: Kesklinna linnaosa Nimi liigisõnaga: Kesklinna linnaosa		Kood: 03QJ Nimi: Müürivahe tn Nimi liigisõnaga: Müürivahe tänav		Kood: 0PEN Nimi: 14 Nimi liigisõnaga: 14		377840298000003QJ00000PEN00000000	
3	390906	3853189	elukondlik hoone	Kood: 37 Nimi: Harju maakond Nimi liigisõnaga: Harju maakond	Kood: 784 Nimi: Tallinna linn Nimi liigisõnaga: Tallinna linn	Kood: 0298 Nimi: Kesklinna linnaosa Nimi liigisõnaga: Kesklinna linnaosa		Kood: 03QJ Nimi: Müürivahe tn Nimi liigisõnaga: Müürivahe tänav		Kood: 0PEN Nimi: 14 Nimi liigisõnaga: 14		377840298000003QJ00000PEN00000000	
4	390906	4808103	elukondlik hoone	Kood: 37 Nimi: Harju maakond Nimi liigisõnaga: Harju maakond	Kood: 784 Nimi: Tallinna Iinn Nimi Iiigisõnaga: Tallinna Iinn	Kood: 0298 Nimi: Kesklinna Iinnaosa Nimi Iiigisõnaga: Kesklinna Iinnaosa		Kood: 03QJ Nimi: Müürivahe tn Nimi liigisõnaga: Müürivahe tänav		Kood: 0PEN Nimi: 14 Nimi liigisõnaga: 14		377840298000003QJ00000PEN00000000	

# Our solution

Demo

## Conclusion

Results vary with different search methods Primitive SQL search

- leaves out results
- includes too many results
- it's much slower (if not optimized)

#### Full text search

- reasonably accurate results
- faster

We are unsure of why Maanteamet's query leaves out some results