

Linguistic Annotation of Old Swedish Texts



Yvonne Adesam

Uppsala, March 2015

MApiR project

- ▶ **M**ethods for the automatic **A**nalysis of **T**ext in digital **H**istorical **R**esources
mapir (*mather, mandr, man* etc.)
OSw: man, person
- ▶ Gerlof Bouma (PI), Yvonne Adesam
- ▶ Investigate and develop natural language processing tools for Old Swedish (1225–1526)
- ▶ 2014–2016/7, funded by the Marcus and Amalia Wallenberg Foundation

Outline

Introduction

Lexical linking

Manual annotation

PoS tagging

MAPiR data

- ▶ The Old Swedish (*fornsvenska*) period starts with latin script manuscripts of law text, ends with (pre-)publication of the Gustav Vasa bible
- ▶ Fornsvenska textbanken, 3.1M tokens
 - ▶ Early (ca 1225-1375)
 - ▶ Late (ca 1375-1526)
- ▶ Svenskt Diplomatarium, 1M tokens (3.5M)

MApiR data

The screenshot displays the MapiR web interface for searching the word "hus". The browser address bar shows the URL: `spraakbanken.gu.se/korp/?mode=old_swedish#?lang=en&stats_reduce=word&cqp=%5B%5D&corpus=fsv-aldrelagor,fsv-aldrelagor,fsv-aldrelagor,fsv-aldrelagor,fsv-aldrelagor`. The page title is "KORP".

The search results are displayed in a table with columns for KWIC (Key Word In Context) and Statistics. The KWIC column shows the word "hus" in various contexts, such as "hus", "hus oc røter oc sum brännir bátrár en [gin] oc en", "Brytir. Y. hus mans oc takir engti bort böte. iij. marker / taki", "ä] [klocn] hus gerneng böte öre firi oc twa öra firi accer en q", "yrklugarp oc hus i sátiá vtán. bispups. orlof oc socnámánná wii", "Nu kan hus niþer falla aff præstins wanrókt; þa a", ". wtan brytae hus hans. dýli mæþ XIII", "andæ. óþaes hus alt. byglæ wp ater. ok bötae", "Bygger hus a humblæ gaarp annars. a. watn wæg æller ky", "Nu far man yfwr mærkia: i akri mans. æller hus sæter a aker annars [före af]", "æng. þa skal han ryblæ þre dyssiar. ok góra eth hus i fira knutær. ok gangi", "winz æy þæt till þa skal hus mætae.", "Il Vm kirkió bol oc hus.", "Nu agho böndær hus a kirkió bol föra. þæt æro siu lagha hus. Stowa.", "Nu agho böndær hus a kirkió bol föra. þæt æro siu lagha hus. Stowa. oc stecara hus.

Other resources

- ▶ Schlyter (1877), *Ordbok till Samlingen af Sweriges Gamla Lagar*, 10 000 entries
- ▶ Söderwall (1884–1973), *Ordbok över svenska medeltids-språket*, 25 000 entries
- ▶ Seed morphology

Historical Material in Computational Linguistics?

- ▶ Low resource / high-variability material
- ▶ Explore methods for such material
- ▶ Understand the material better by looking at what methods work
- ▶ Results may carry over to other types of material, e.g. social media

Challenges of OSw text

- ▶ Lexicon: *aptanbakka, bakvapi*
- ▶ Morphology, syntax: *Nu kan kirkia brinna ælla stulin uarþa [at] ipnum durum*
- ▶ Lack of standardization:
 - ▶ spelling: *bokstaffua, bokstaffwa, bokstafwa, bokstaua, bokstawa, bogstaffwa*
 - ▶ segmentation: Sentence units marked by period, slash, comma, capital, or not at all
- ▶ Lack of resources: No annotated data, no complete descriptions, no native speakers

Our approach to NLP for OSw

- ▶ Linguistically informed
- ▶ Linking
 - ▶ lexicon to lexicon
 - ▶ text to lexicon
 - ▶ between spelling variants
 - ▶ between text editions (e.g. bible for annotation projection)
- ▶ Manual annotation for training/evaluation data
- ▶ Annotation projection (from Icelandic, modern Swedish)

Linking to lemmata

- ▶ *bokstaffwa bokstaff bokstawom bokstaf
bogstaffwa bokstaua bokstawa bokstafwa
bokstaffwom bokstaffua bokstawin
bokstaffwane bokstaffwinor(=m?) bokstafuom
bokstawane*
- ▶ How do we recognize these as occurrences
(inflections + spelling variants) of **bokstaver**
(Söderwall)

Linking to lemmata 2

- ▶ Using lexicon entries as lemmata
 - ▶ language of the period
 - ▶ information about PoS, morphology etc
 - ▶ Definition and description

A simple algorithm for linking

- ▶ List of lemmata + occurrences from dictionary
- ▶ b o g s t a f f w a _
b o k s t a _ _ v e r
(iterated LD alignment: Wieling et al., 2009)
- ▶ Count frequent character group substitutions, turn into rules, weigh according to frequency
- ▶ Link a token to the lemma that requires the cheapest substitutions

A simple algorithm for linking

- ▶ On *Abota*:
 - direct lookup 2/10 correct on first guess
 - simple linking 5/10 on first guess
- ▶ Split inflection from spelling variation?
 - use explicitly coded human knowledge (computational morphology)?

Mozilla Firefox
http://dem...agar/Ogl-A x
demo.spraakdata.gu.se/fsvreader/svndict/showtext/big/aldre_lagar/Ogl-A Google

Meny Tillbaka

XIII.

Nu sitær bonde kuar mæþ tiunde um ar. böte firi þrea öra. ok ut tiundan præstinum. ælla þöm sum tiundan a. sua firi annat ok sua firi þriþia. Sitær kuar um þry böte biskupe þrea markær. Nu före bonde egh tiunda til firi paska. þa halde þöghine hans rezskap [up] utan kalle han firi kirkju dyr æpte paska. ok uiti þæt mæþ eþe fiughurtan manna at han sat mæþ luui kuar. ælla ok at han förþe han. ælla han böþ han ut. ok þe wltu sialuir at þe uilldu egh uipær taka. orka egh eþe böte þrea öra. æn han sat minna kuar æn um þry ar. sat han um {þru} [þry] ar böte þre markær sum för uar saght. §. 1. Nu ær biskupær skyldughær firi tiunda sin. wighia krismu ok klærka. kalk mæssu klæþe: huart þriþia ar {tik} [til] sokn koma tuægga natta gingærþ af præste taka. mæþ mannum tolf ok sialuær han þrattande þa a han folk færma ok næmdir sea. §. 2. þa biskupær uill næmd sea. þa skal han fa manaþa buþ firi sik. præstær a næmd gæra sanna mæn i næmd sætia ok egh uipærdelu mæn. ok egh soknara. ok egh þæn i epum in stop. ok þe skulu sitia i næmdinne sum þe sighia ia uipær sum uipærdulu mænnini æru. sipan a soknarin eþa til næmd bæra. þa skal næmdin sitia sik ensamin ok talas uipær. alla þa eþa sum hon uær þa uarin uarþi. alla þa eþa* hon fælle þa uarin fælde. utan næmdin uili nakuara eþa undi kunungx dom læggia. ælla laghmanzs. þa biþin þe eþa þær til þe uarþa fællde ælla uarþe. §. 3. all laghkallaþ mal aghu a biskups næmd kuma ekki huma mal. Alli eþa aghu a biskups næmd kuma. þæssin mal aghu a biskups næmd koma. hor. menepa.

demo.spraakdata.gu.se/fsvreader/lexeasy/þrattande-þrattande-þratta-þrætta-þrättan-þrættan

þrattande 5 träffar. (þrattande/þratta/þrætta /þrättan/þrættan)

- þrattande nl
 - **Söderwall nl** trettonde. " thässe fäm konunga hafdho thiánt tolf aar mz skat vnder kodorlaamor konung, oc thrättanda aarith gingo the aff honum " MB 1: 180 . Bir 4: 23 . VKR 69 . Lg 554 .
 - **Söderwalls Supplement nl** trettonde. " j the kirkio standa tolf stola oc thrättande war then som gudh siäff sat a Prosadikter (Karl M) 251. ib. ther waro tolf sänga . . . oc thrättande i midhio " ib 255 . - þrattande dagher iula, trettonde dag jul, trettondedag(en). trydie torsdagen nest efter trettendhe ... [more](#)
- þratta vb nn
 - **Söderwall vb L**.
 - **Söderwall vb** envisas, anstränga sig,

Manual annotation

- ▶ Mainly following guidelines for Old Norwegian from Menotec project (Haugen & Øverland 2014)
- ▶ Lemmata, morphological features, syntactic structures using the PROIEL scheme (Eckhof & Haug 2010) and annotation environment – also used for 16(!) other historic language annotated corpora.
- ▶ Lemmata: Söderwall entries
- ▶ Status: almost 20k tokens annotated

Data

- ▶ Östgötalagen (A Holm B50), ca 1290/1350
18'269
- ▶ Ä. Västgötalagen (A Holm B59), ca 1220/1280
491
- ▶ Abota (Klemming), ca 1448
541
- ▶ Pentateuchparafr. (B Holm A1), ca 1330/1526
469

Data

- ▶ Caveats:
 - ▶ no punctuation
 - ▶ manual sentence splitting
 - ▶ manual tokenization

Data

- ▶ Caveats:
 - ▶ no punctuation
 - ▶ manual sentence splitting
 - ▶ manual tokenization

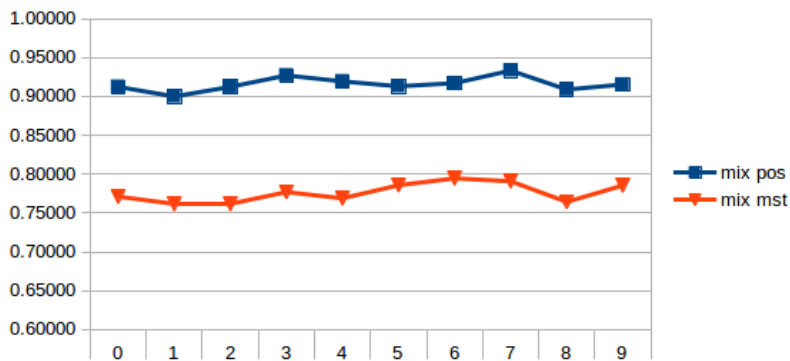
*Uerder mæper .i. kyrkiu dræpin þet ær nipings værk.
þa er kyrkia al vuighz.*

If a man is killed in church, this is infamy, then the whole church is unhallowed.

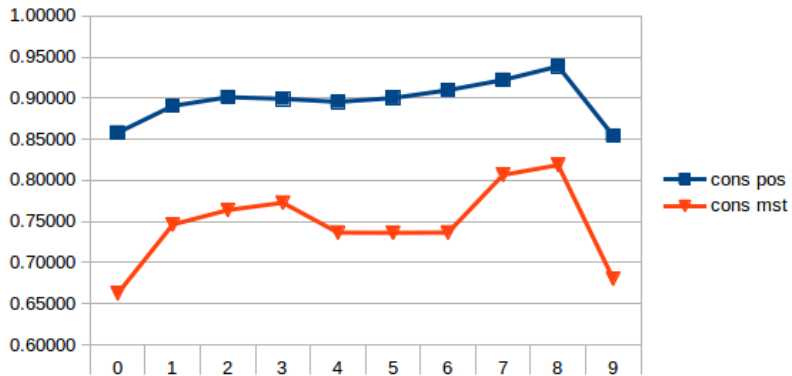
Östgöta crossvalidation

- ▶ Stagger (Östling 2013)
- ▶ Mixed 10-fold
- ▶ Consecutive 10-fold
- ▶ 4-fold by book (*balk*, 3500-5000 tokens)

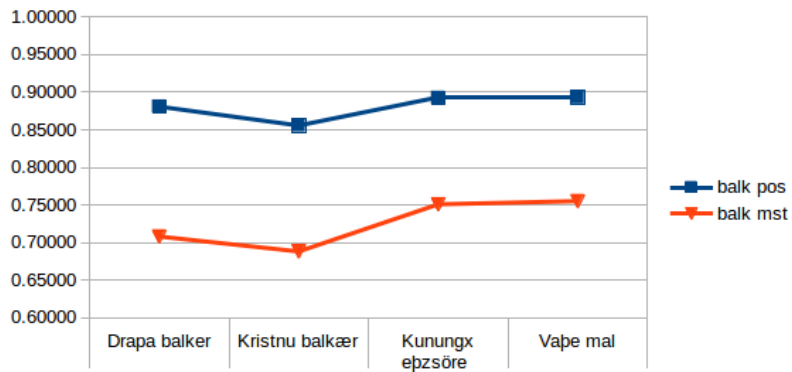
Data splitting



Data splitting



Data splitting



Lemmas and spelling

- ▶ Using lemmas as extra clues in tagging
- ▶ Simple spelling normalization

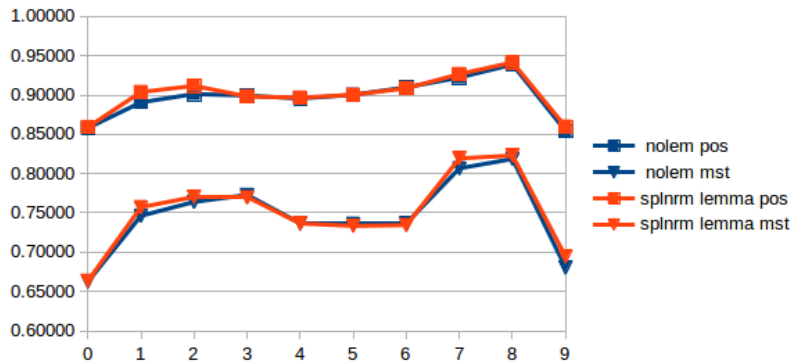
aa \rightarrow a,

ch/c/q \rightarrow k,

th \rightarrow t,

etc.

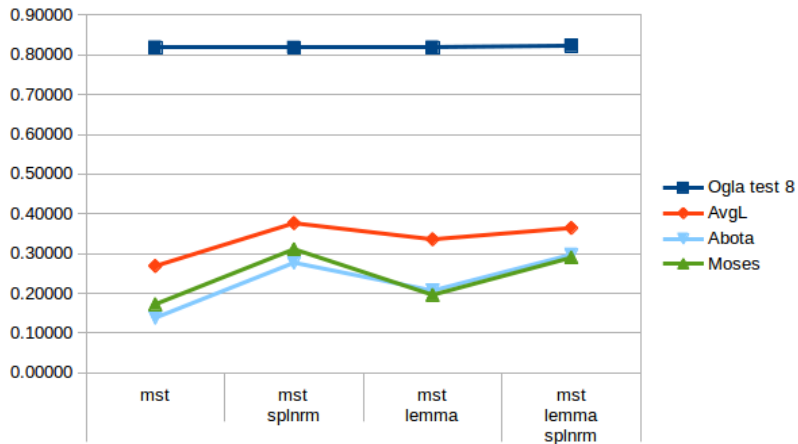
Lemmas and spelling



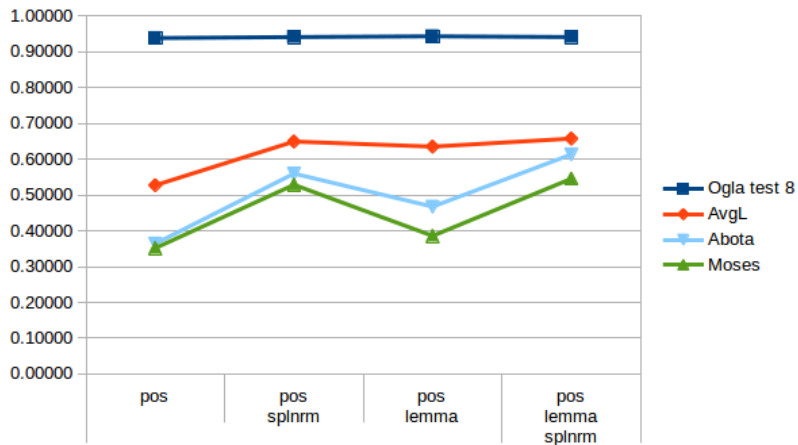
Other texts

- ▶ Östgöotalagen homogenous
- ▶ Testing on other texts?
- ▶ Training on other texts?

Other texts



Other texts



OOV

	Östgl	ÄVästgl	Abota	Mose
Tokens	7%	66%	78%	79%
Types	22%	75%	85%	86%

With spelling normalization:

	Östgl	ÄVästgl	Abota	Mose
Tokens	7%	51%	60%	54%
Types	21%	66%	76%	72%

Lemmas and spelling

Improvements

- ▶ ÄVästgl
 - ▶ spelling: conj, verbs
 - ▶ lemma: verbs, conj
- ▶ Abota
 - ▶ spelling: conj, verbs, adverbs, nouns, prep
 - ▶ lemma: verbs, adverbs, conj
- ▶ Mose
 - ▶ spelling: conj, verbs
 - ▶ lemma: verbs

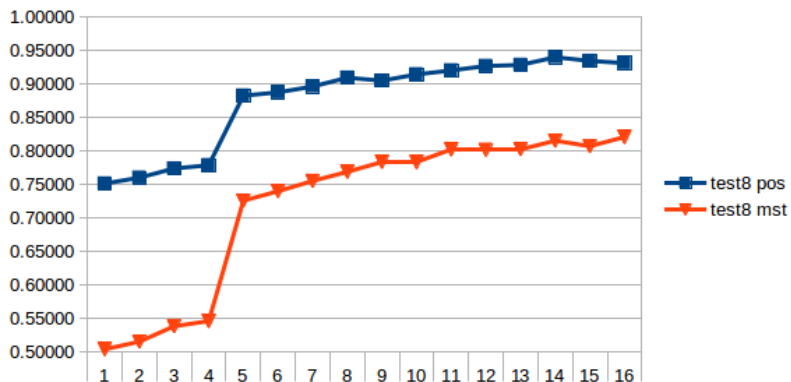
Other texts

- ▶ Training 5k Östgöta + 500 other text
- ▶ Inconclusive results
- ▶ Adding another text sometimes beneficial (compared to just using a lot more Östgöta)
- ▶ A lot more Östgöta always better when applying spelling normalization

Training size

- ▶ How much data is needed?
- ▶ Increase training size in steps

Training size



Conclusions

- ▶ Explore linguistically informed methods for annotation of OSw
- ▶ PoS experiments on manually annotated data
- ▶ Spelling normalization and lemmas improve results (autom. lemmas?)
- ▶ Fairly small amounts of annotated data necessary
- ▶ Large variation between texts