

A Thesaurus-Based Sentiment Lexicon for Danish – The Danish Sentiment Lexicon

Sanni Nimb¹, Sussi Olsen², Bolette S. Pedersen², Thomas Troelsgård¹

¹The Society for Danish Language and Literature & ²University of Copenhagen
sn@dsl.dk, saolsen@hum.ku.dk, bspedersen@hum.ku.dk, tt@dsl.dk

The Danish Sentiment Lexicon characteristics

- 13,859 Danish polarity lemmas
- -3 to +3
- Includes morphological information
- Freely available at <https://github.com/dsl/danish-sentiment-lexicon> (licence CC-BY-SA 4.0 International)

Main idea

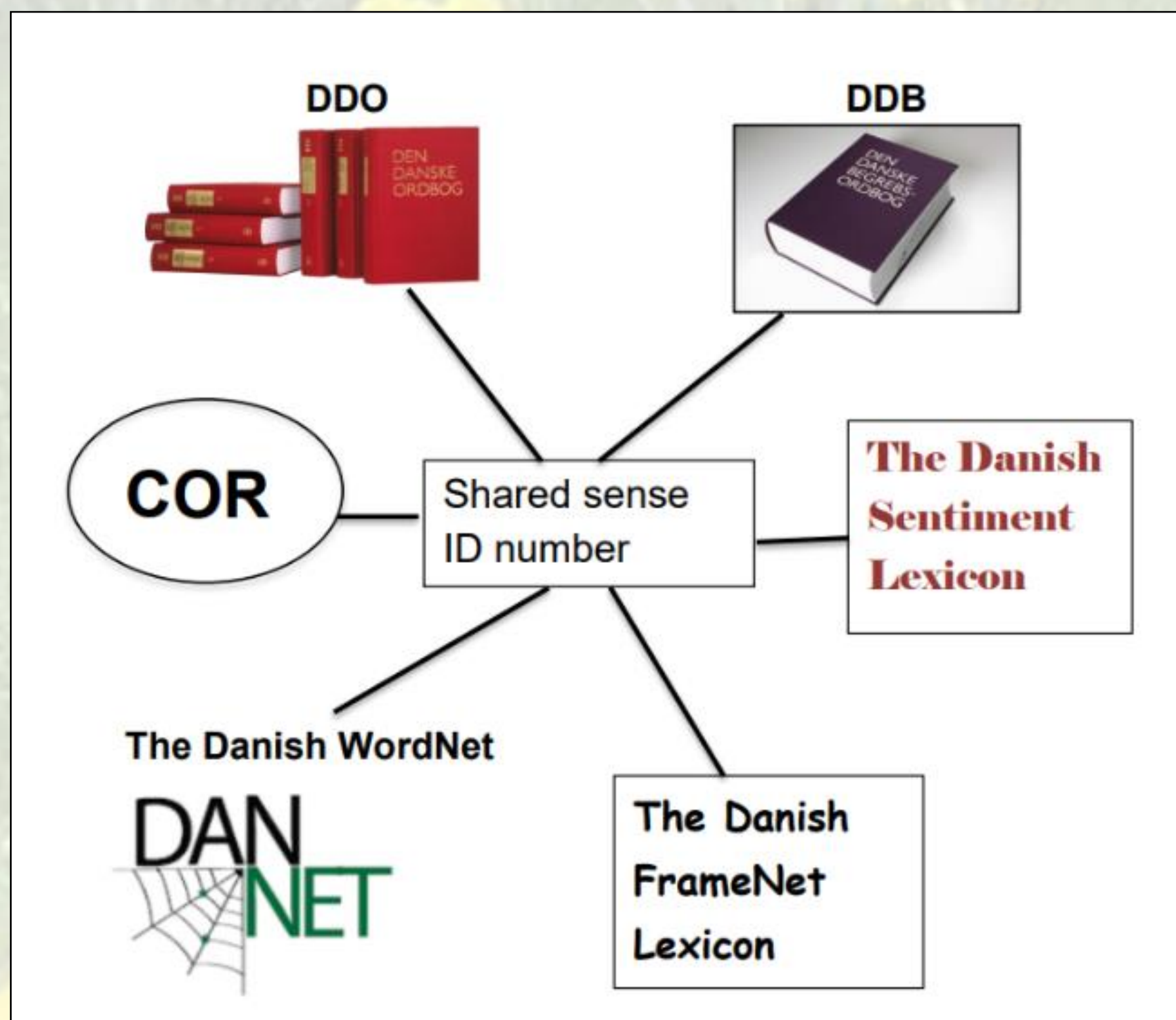
- Easy access to comprehensive vocabulary in thesaurus → much higher lexical coverage than existing Danish sentiment lexicons
- Identify positive and negative thesaurus section titles as the starting point
- Calibrate the degree of polarity within sets of words denoting very similar concepts

Linked data

- Based on 17,883 polarity annotated senses linked to other Danish lexical resources
- Differentiates it from other sentiment lexicons for Danish
- Allows for future experiments where sentiment is combined with other types of information from lexicons and corpora

Linked Data

- Four computational lexicons linked to a monolingual dictionary DDO and to the thesaurus DDB via the sense ID numbers of DDO
- The WordNet linked to Princeton WordNet

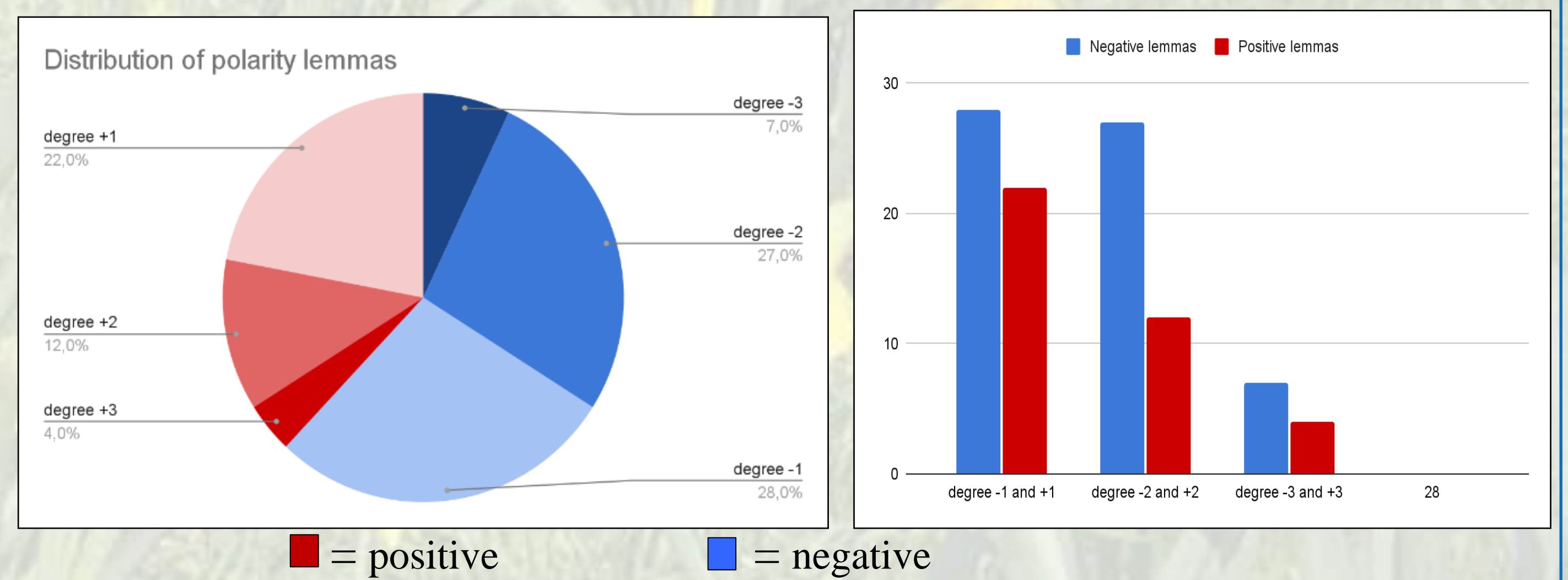


Annotation of the dataset

	Task	The judgement is based on
Step 1	sense polarity: 0, positive or negative?	<ul style="list-style-type: none"> • AFINN • Comparison of synonyms and near synonyms in the DDB thesaurus • Information in DDO
Step 2	degree of sense polarity: -3, -2, -1, +1, +2, +3?	<ul style="list-style-type: none"> • AFINN • Comparison of synonyms and near synonyms in the DDB thesaurus • Information in DDO
Step 3	harmonizing values of identical lemma senses in dataset	<ul style="list-style-type: none"> • Lemma / sense represented more than once in the dataset? (due to multiple representations in thesaurus)
Step 4	deciding upon polarity at lemma level	<ul style="list-style-type: none"> • Conflicting polarities of lemma senses? • Rare sense to be ignored? • Or lemma to be left out?
Step 5	validation of data	<ul style="list-style-type: none"> • 1/3 of annotated data • Comparison of all lemmas with same high degree (e.g. all +3 lemmas)

Results and conclusions

- Clear predominance of negative lemmas
62% negative polarity
38% positive polarity
- ~ same distribution as in Swedish SenSALDO (6,386 lemmas, also thesaurus-based, see Rouces et al., 2018a and b)
- Devitt & Ahmad (2013): seems general for sentiment lexicons; in texts the relationship is thought to be reversed.



→ When sentiment analysis methods are based on large lexicons (like ours) the variation of lemmas that are recognized is higher – and these are more likely to be negative than positive → overweight of negative polarity
→ When methods are based on smaller lexicons, it is more likely that the neglected lemmas are negative than positive → underrepresentation of negative polarity

- We manually annotated a literary piece of text containing many polarity conveying lemmas (Pedersen et al. 2021)
→ only ad hoc composita not covered → we should include morphemes with high polarity+ module for automatic splitting of unknown composita

- There might be undiscovered polarity lemmas in thesaurus/DDO
 - Hidden in neutral thesaurus section (e.g. 'Man', 'Woman'). Should be added.

- We might consider including neutral lemmas (like SenSALDO)

- Plan: transfer of polarity information at sense level (not part of the open release)
 - to DDO → new presentations, e.g. of synonyms
 - to DanNet
 - to new formal semantic lexicon COR-S (Pedersen et al. LREC 2022)

Creation of the dataset

• 24% of the 888 sections judged by two lexicographers to contain polarity words based on the section name. Clear predominance of negative sections:

- 57% negative (e.g. sections *Unimportant*; *Sadness*)
- 37% positive (e.g. sections *Important*; *Admire*; *Friendship*)
- 6% both negative/positive (e.g. sections *Reputation*; *Protest*, *Rebellion*).
- → in total 25,000 senses (not including MWU).

• Senses from identified sections extracted from DDO and as default assigned the section polarity value

• Combined with information from other resources

- DDO (definition, frequency, usage (e.g. *derogatory*; *offensive*))
- Values from AFINN (Danish, ~3000 words, -5 to +5, see Nielsen, 2018)
- Keyword and synonyms/near synonyms in thesaurus sections

7	9	19	Uppgive, give arkain pa	vero	3	abandonnere	26114378	abandonnere	oppgive, give slip paopgive	n	-2	0	give arkain pa, give	
7	9	7	Tvivlrtdig, forvirret	adj	3	iben	21098037	iben	(endnu) ikke afklaret, afgjort	0	-2	2	uvist, ukendt, uklar, u	
7	9	20	Handic, agere	subst	3	aning	21098099	aning	begrndelse, indledning, beg	0	2	0	handling, forehavend	
7	9	3	Vere nedt til	adj	3	absolut	21000113	absolut	som er ubetinget og uafhaer p	0	-2	2	0	nodvendig, uundgael
7	9	6	Beslutsum	subst	3	absolutthed	21000116	absolutthed	viljen til at gennemfore nog p	0	2	1	beslutsumhed, deterr	
7	9	6	Beslutsum	subst	3	absoluttet	21000119	absoluttet	person som glir ind for absp p	0	2	0	en handlingens mand	
7	9	41	Vigtig	subst	3	accent	21000162	accent	(ofa)hovedvaegt; eftertryk;	0	2	0	prioritering, vaegtning	
7	9	41	Vigtig	verb	3	accentuere	21000167	accentuere	fremhaeve; understrege/frem p	2	0	0	fremhaeve, pointere,	

Inter-annotator agreement on sample

- Two lexicographers annotated 400 words also present in AFINN.
- Negative and positive values, but not degree → Cohen's Kappa: 0.83