

EXPLORING DIFFERENT WORD EMBEDDINGS FOR UNSUPERVISED PART-OF-SPEECH TAGGING

MSc dissertation

Felicia Liu

November 14, 2016

- Introduction to thesis topic
- Research objective
- Methods & Approach
- Results

INTRODUCTION TO RESEARCH TOPIC

Exploring word embeddings for unsupervised part-of-speech tagging.

POS tagging Assigning labels to individual words from a given text that denote its syntactical function

Unsupervised Only use unlabeled data

Word embeddings Vector representations of words, trained to encode a word's features

Use Gaussian Hidden Markov Model for POS tagging.

RESEARCH OBJECTIVE

- Use word embeddings to improve POS tagging.
- Encode relevant information in these embeddings, such as syntax or morphology.

How do we create such word embeddings?

- Small context size window
- Character-level word embeddings

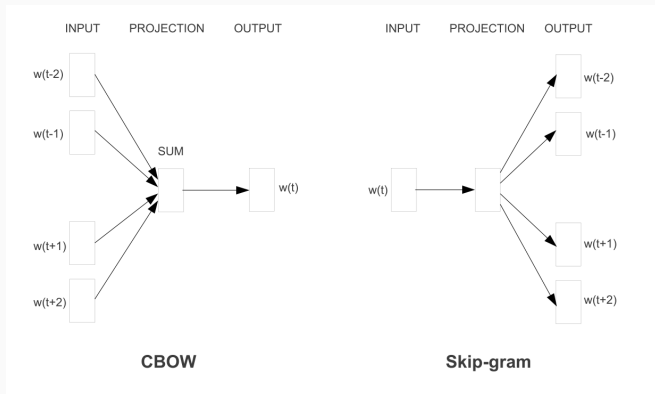
METHODS & APPROACH

- Use one data set to train all embeddings.
- Train embeddings with vector sizes 20, 50, 100, and 200, vary other parameters
- Use V-measure to assess performance with gold-standard labeled text.

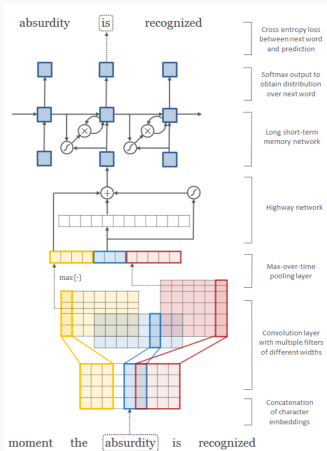
Word embeddings used in this work include

- SENNA (Collobert et al., 2011)
- GloVe (Pennington et al., 2014)
- word2vec (Mikolov et al., 2013)
- Structured word2vec (Ling et al., 2015)
- Character-level embeddings (Kim et al., 2016)

WORD2VEC MODEL

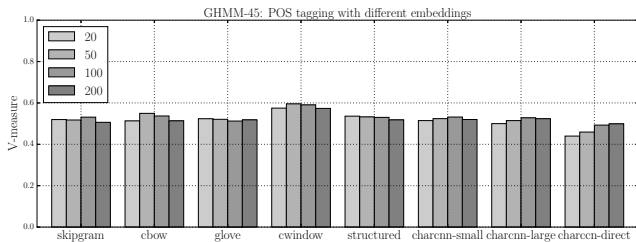


CHARACTER-LEVEL EMBEDDINGS MODEL

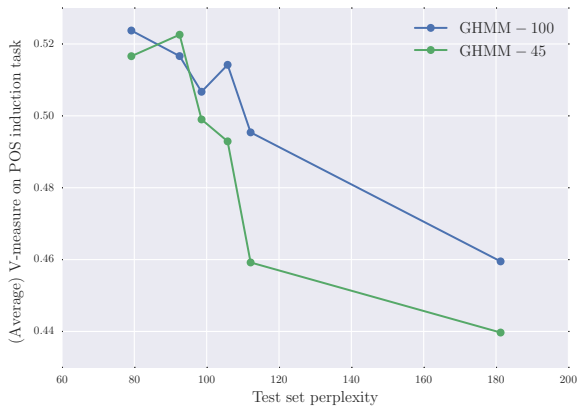


- Varying the level of tokenisation of input data
- Reducing dimensionality of embeddings obtained through the character-level model
- Combining word- and character-level embeddings

RESULTS ON POS TAGGING



INFLUENCE OF PERPLEXITY ON RESULTS



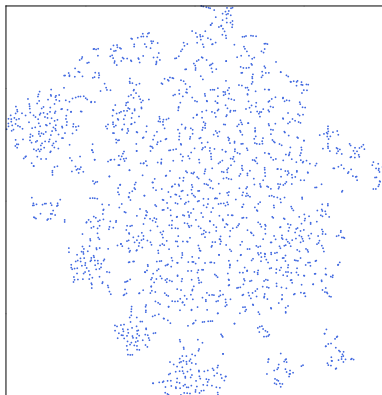
NEAREST NEIGHBOURS 1

word	skip-100	CBOW-50	GloVe-20	cwin-50	struct-20	SENNA	Large-100	Small-100
person	woman	woman	true	woman	woman	letter	peterson	peterson
	persons	child	man	child	firm	case	persons	pension
	anyone	settlor	result	man	child	honor	pearson	persons
	patient	spouse	real	nation	hobby	ages	emerson	pearson
	victim	persons	woman	girl	demon	problem	patterson	poisson
king	kings	hussa	prince	prince	queen	lord	ping	ping
	queen	theodric	alexander	queen	lord	queen	ming	qing
	harthacanute	kings	charles	captain	lady	emperor	qing	kind
	mordha	sillok	edward	bishop	prince	titles	kind	bing
	monarch	culen	henry	lord	grace	fighting	ring	kong
france	spain	belgium	rome	spain	cuba	santa	franc	trance
	italy	italy	germany	belgium	luxembourg	composition	franco	franc
	belgium	spain	spain	italy	guatemala	germany	frances	franco
	germany	bordeaux	portugal	austria	portugal	italy	franca	ordnance
	luxembourg	luxembourg	japan	poland	peru	sweden	frances	franca
reddish	grayish	grayish	lime	yellowish	whitish	violet	resisted	yiddish
	greyish	greenish	honey	grayish	wavy	territorial	yiddish	swedish
	yellowish	blackish	trout	bluish	feathered	academia	revised	rush
	greenish	purplish	yellow	mottled	coppery	aggregate	registered	dish
	bluish	irides	fat	red	bulbous	tore	swedish	irish
richard	robert	robert	robert	william	william	robert	richards	orchard
	walter	philip	george	harold	harold	peter	richardson	richards
	francis	william	francis	robert	stephen	reportedly	orchard	richardson
	hugh	ralph	james	charles	albert	david	richland	richland
	arthur	john	thomas	hugh	john	william	archaic	richmond

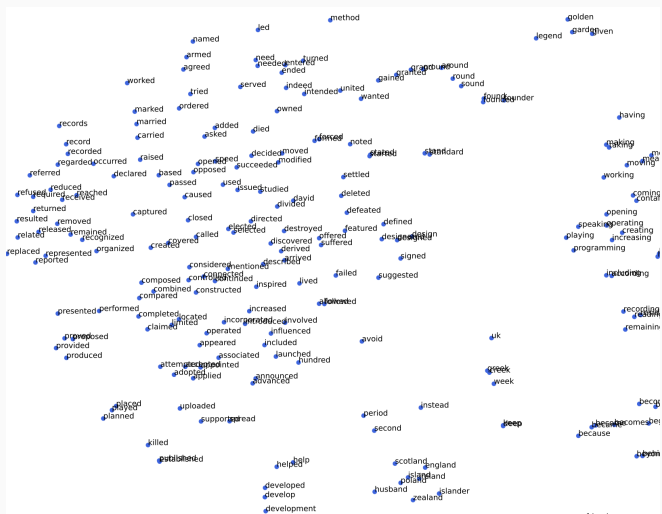
NEAREST NEIGHBOURS 2

word	skip-100	CBOW-50	GloVe-20	ewin-50	struct-20	SENNA	Large-100	Small-100
inconsiderable	denominate policyholders significant extortionate substantial	sphenic discounting surprising denormal staggering	endear materialise disappoint congenial distasteful	obligatory significant unimpeachable unaffordable explicit	destabilizing satisfactory foolproof scientifically realigning		considerable inconquerable inconsolable insufferable imponderable	considerable indecomposable inconquerable incommensurable unconscionable
unsteadiness	neuralgic schizotypy insensibility vestibulo unnaturalness	equilibrioception neuralgic persecutory schizotypy diffractive	zelotes christopher wck ranko hypernatremia	decompensated dyserythropeitic microcornea expressivity morphoea	thymoma mechanical hemiplegic hyperglycemic indirectness		unreadiness steadiness untidiness uneasiness steadfastness	unreadiness uneasiness steadiness untidiness sturdiness
commenting	commented gloating raved remarked joked	commented criticizing insisting remarking discussing	insisting sells binds insists concentrating enabled	insisting insists focusing commenting speculating insisted	insisted insisting commented speculating insists	hackers possessed corvette cyborg jtdirl	committing competing commanding coming connecting	committing commemorating commanding coming competing
comment	comments remarks reply remark quip	comments reply remark remarks quip	request notice finding account permission	remark report notice statement commentary	consensus slump shame ban commentary	lyon marco nan thebes orchestras	commencement commitment competent clement comments	commitment commencement component competent clement
unaffected	affected obscured disturbed untouched evident	obscured disturbed hindered affected compromised	disrupted outdated motivated unreliable unstable	affected disrupted regulated damaged overwhelmed	affected disrupted obscured enforced regulated	micronesia baptist brett apparatus trenton	affected unwanted uninhabited unidentified unidentified	infected affected inflicted effected unidentified
affect	impair contribute depend affecting affects	contribute impair reflect alter relate	reject observe recognize deny arise	reflect disrupt eclude involve satisfy	utilize involve activate utilise overwhelm	ngo paula cabin collier mentally	effect affects affected affecting affection	effect affects affected affecting defect

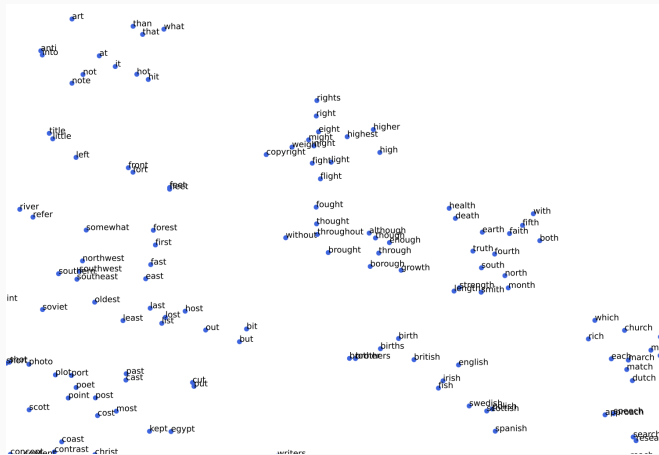
T-SNE APPLIED TO CHARACTER-LEVEL EMBEDDINGS



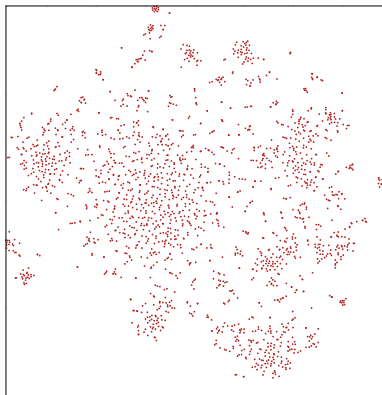
CHARACTER-LEVEL EMBEDDINGS CLUSTERS 1



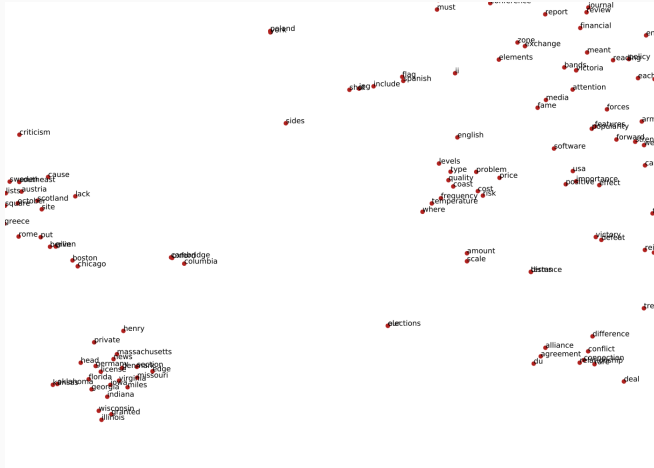
CHARACTER-LEVEL EMBEDDING CLUSTERS 2



T-SNE APPLIED TO VISUALISE WORD-LEVEL EMBEDDINGS



WORD-LEVEL EMBEDDING SPACE CLUSTER



CONCLUSION

- Structured word embeddings do better than unstructured.
- Continuous bag-of-words (CBOW) performs better than skipgram.
- Character-level embeddings do not show competitive performance when used alone, can improve performance when combined.
- Perplexity of language model is correlated with embedding performance on POS tagging task.

QUESTIONS?