# DISTRIBUTIONAL SEMANTICS AND COMPOSITIONALITY

Corina Dima

25rd April 2019

#### REVIEW

use the co-occurrence information in large corpora to create distributional representations for words

	iced	(to) drink	owner	in
cappuccino	6	2	0	3
espresso	1	1	0	4
cat	0	1	4	3
latte	6	5	0	4
leaf	0	0	0	5





$$f\left(0.3\,0.1\,0.7\,u\,,\,0.5\,0.9\,0.1\,v\right) = ??\,??\,p$$

What f makes p most similar to w?

 $p = g(\mathcal{W}[u \odot u'; v \odot v''] + b)$ where  $p, u, u', v, v'', b \in \mathbb{R}^{n}; W \in \mathbb{R}^{n \times 2n}; g = tanh$ 

 $\begin{aligned} & \underset{p \in \mathcal{W}g(\mathcal{W}_{_{1}}[u; v] + b_{_{1}}; \mathcal{W}_{_{2}}[u; v] + b_{_{2}}; ...; \\ & \mathcal{W}_{_{k}}[u; v] + b_{_{k}}) + b \end{aligned} \\ & \text{where } p, u, v, b, b_{_{i}} \in \mathbb{R}^{n}; \mathcal{W}_{_{i}} \in \mathbb{R}^{n \times 2n}; \mathcal{W} \in \mathbb{R}^{n \times kn}; g = relu \end{aligned}$ 

#### COMPOSITIONALITY

- Composition models for distributional semantics can extend the vector spaces
- they learn how to build representations for complex words (e.g. 'apple tree') and phrases (e.g. 'black car') from the representations of individual words.

- ► when using a search engine to look for job ads
  - ► programmer
  - ► software developer



- ► when using a search engine to look for job ads
  - ► programmer
  - ► software developer



- ► when using a search engine to look for job ads
  - ► programmer
  - ► software developer

- when shopping online
  - ► blue jacket
  - ► dark-blue coat





. . . . . . . . . . . . . . . . . . .

. . . . . . . . . . .

#### Handpuppentheateraufführung

#### Handpuppentheateraufführung Hand | Puppen | Theater | Aufführung

Handpuppentheateraufführung Hand | Puppen | Theater | Aufführung hand | doll-s | theater | performance





















2,019,439 35,921 408,726 91,268 Hand | Puppe(n) | Theater | Aufführung













Handpuppentheater

 2,019,439
 35,921
 408,726
 91,268

 Hand
 Puppe(n)
 Theater
 Aufführung





Handpuppentheater

 2,019,439
 35,921
 408,726
 91,268

 Hand
 Puppe(n)
 Theater
 Aufführung



## 66

"The meaning of the whole is a function of the meaning of the parts and their mode of combination"

Frege's Principle of Compositionality

. . . . . . . . . . . . .

. . . . . . . . .

- use distributional representations as a proxy for the meaning of individual words - "meaning for the parts"
  - ► any type of vector representations for words

- use distributional representations as a proxy for the meaning of individual words - "meaning for the parts"
  - any type of vector representations for words
- Iook for a "mode of combination"

- use distributional representations as a proxy for the meaning of individual words - "meaning for the parts"
  - > any type of vector representations for words
- Iook for a "mode of combination"
- ► in the context of vector space models of language
  - find a composition function f that fulfils certain requirements

. . . . . . . . . . . . . .

. . . . . . . . .

*f* takes as input (two) *d*-dimensional distributional
 representations of the constituents (learned from a corpus)

*f* takes as input (two) *d*-dimensional distributional
 representations of the constituents (learned from a corpus)

 $f(u^{corpus}, v^{corpus})$ 

*f* takes as input (two) *d*-dimensional distributional
 representations of the constituents (learned from a corpus)

 $f(u^{corpus}, v^{corpus})$ 

f produces as output another *d*-dimensional vector, which is a combination of the input vectors

*f* takes as input (two) *d*-dimensional distributional
 representations of the constituents (learned from a corpus)

 $f(u^{corpus}, v^{corpus})$ 

*f* produces as output another *d*-dimensional vector, which is a combination of the input vectors

$$p^{composed} = f(u^{corpus}, v^{corpus})$$



#### Hand 'hand'





Puppe 'doll'

0.1	v
0.3	V
0.2	



Handpuppe 'hand puppet'

0.4 0.3 *p* 0.2



#### Hand 'hand'





Puppe 'doll'

0.1	1,
0.3	V
0.2	



#### Handpuppe 'hand puppet'

0.4	
0.3	$\widetilde{p}$
0.2	

original (learned from the corpus)



Hand 'hand'





Puppe 'doll'





#### Handpuppe 'hand puppet'

0.4	
0.3	$\widetilde{p}$
0.2	

???

p

original (learned from the corpus)

$$f(\begin{array}{cccc} 0.2 \\ 0.1 \\ 0.4 \end{array}, \begin{array}{cccc} 0.1 \\ 0.3 \\ 0.2 \end{array}, \begin{array}{cccc} 0.1 \\ 0.3 \\ 0.2 \end{array} =$$



Hand 'hand'





Puppe 'doll'

 $\begin{array}{c} 0.1 \\ 0.3 \\ 0.2 \end{array}^{\nu} ) = \\ \end{array}$ 



0.2 0.1 0.4

f(



#### Handpuppe 'hand puppet'

0.4	
0.3	p
0.2	

?

p

original (learned from the corpus)





Hand 'hand'





Puppe 'doll'



Handpuppe 'hand puppet'

0.4	
0.3	Ĩ
0.2	

original (learned from the corpus)

$$f(\begin{array}{cccc} 0.2 \\ 0.1 \\ 0.4 \end{array}, \begin{array}{cccc} 0.1 \\ 0.3 \\ 0.2 \end{array}, \begin{array}{ccccc} v \\ 0 \end{array} = \begin{array}{ccccc} ? \\ ? \\ ? \end{array} p$$

composed

$$\operatorname{argmax} \frac{p \cdot \tilde{p}}{\|p\|_2 \|\tilde{p}\|_2}$$

#### ADDITIVE COMPOSITION - MITCHELL & LAPATA 2010



#### ADDITIVE COMPOSITION - MITCHELL & LAPATA 2010



### ADDITIVE COMPOSITION - MITCHELL & LAPATA 2010



#### "You might just as well say that 'I see what I eat' is the same thing as 'I eat what I see'!"

(The Mad Hatter in Lewis Caroll, Alice's Adventures in Wonderland)

### WEIGHTED ADDITIVE COMPOSITION - MITCHELL & LAPATA 2010


### WEIGHTED ADDITIVE COMPOSITION - MITCHELL & LAPATA 2010



## WEIGHTED ADDITIVE COMPOSITION - MITCHELL & LAPATA 2010



#### **NEIGHBOURS OF ADDITIVE COMPOSITION**

initiation_rite:1	sewing_machine:1	poverty_trap:1	dipole_antenna:1
initiation_rite 1.00000	sewing_machine 1.00000	poverty_trap 1.00000	dipole_antenna 1.00000
initiation_rite_c 0.48302	sewing_machine_c 0.59397	poverty_trap_c 0.48964	dipole_antenna_c 0.58292
initiation 0.41953	sewing 0.57117	benefit_system 0.46336	dipole 0.51744
kickback_scheme 0.39792	quilting 0.48184	welfare_system 0.41871	antenna 0.46497
prison_program 0.38725	machine_stitch 0.47350	marriage_penalty 0.40658	radiation_pattern 0.44861
rite 0.38340	machine_operator 0.46966	means_test 0.38529	dish_antenna 0.44799
ritual_killing 0.37734	embroidery 0.46061	rat_race 0.38416	throwing_knife 0.42321
rescue_vehicle:100	reaction_time:100	advance_copy:100	phase_space:100
rescue_vehicle 1.00000	reaction_time 1.00000	advance_copy 1.00000	phase_space 1.00000
rescue_equipment 0.49602	working_memory 0.51014	review_copy 0.59946	space_time 0.44398
rescue_squad 0.45921	response_time 0.47491	draft_copy 0.56645	phase_transition 0.42528
rescue_personnel 0.45826	learning_ability 0.46824	paperback_copy 0.51407	ground_state 0.41160
rescue_crew 0.45364	processing_time 0.44939	sneak_preview 0.47378	field_theory 0.41095
rescue_vessel 0.42381	cycle_time 0.44638	galley_proof 0.41045	quantum 0.40099
rescue_vehicle_c 0.30007	reaction_time_c 0.27750	advance_copy_c 0.27408	phase_space_c 0.26103
paper_market:1000	vacuum_aspiration:1000	showpiece_event:1000	cart_track:1000
paper_market 1.00000	vacuum_aspiration 1.00000	showpiece_event 1.00000	cart_track 1.00000
property_lending 0.63559	self_stimulation 0.46542	surprise_team 0.55356	dirt_path 0.58237
property_share 0.50744	crystal_gazing 0.45362	budget_session 0.54997	go_cart 0.56240
aircraft_market 0.50299	garden_cart 0.45026	weekend_election 0.53930	dirt_track 0.52775
property_sector 0.48177	exchange_transfusion 0.44189	rating_period 0.53279	logging_road 0.49539
air_lane 0.47996	motorized_wheelchair 0.43895	year_area 0.52381	gravel_road 0.47347
paper_market_c 0.16682	vacuum_aspiration_c 0.07520	showpiece_event_c 0.28813	cart_track_c 0.23656

#### NEIGHBOURS OF WEIGHTED ADDITIVE COMPOSITION

election_campaign:1	decree_nisi:1	dung_beetle:1	drinking_alcohol:1
election_campaign 1.00000	decree_nisi 1.00000	dung_beetle 1.00000	drinking_alcohol 1.00000
election_campaign_c 0.70908	decree_nisi_c 0.51042	dung_beetle_c 0.53302	drinking_alcohol_c 0.64101
election 0.65634	decree 0.45393	stag_beetle 0.45368	smoking 0.61267
fall_election 0.58337	annulment 0.45290	beetle 0.42449	drinking 0.60360
election_cycle 0.57361	summary_judgement 0.44723	pill_bug 0.41598	alcohol_consumption 0.57452
campaign_period 0.57320	committal_order 0.39749	dung 0.39514	binge_drinking 0.54170
campaign_team 0.55469	divorce 0.39232	water_strider 0.38161	alcohol_use 0.53805
rescue_vehicle:100	application_procedure:100	$music\_market:100$	phase_space:100
rescue_vehicle 1.00000	application_procedure 1.00000	music_market 1.00000	phase_space 1.00000
rescue_equipment 0.49602	application_process 0.67551	music_world 0.58263	space_time 0.44398
rescue_squad 0.45921	application_form 0.59632	music_chart 0.54964	phase_transition 0.42528
rescue_personnel 0.45826	job_application 0.54862	pop_music 0.52977	ground_state 0.41160
rescue_crew 0.45364	grant_application 0.50493	music_industry 0.51918	field_theory 0.41095
rescue_vessel 0.42381	visa_application 0.46575	music_sale 0.51078	quantum 0.40099
rescue_vehicle_c 0.29996	application_procedure_c 0.30313	3 music_market_c 0.34352	phase_space_c 0.26121
shoo_fly:1000	tobacco_market:1000	paper_wasp:1000	sports_vehicle:1000
shoo_fly 1.00000	tobacco_market 1.00000	paper_wasp 1.00000	sports_vehicle 1.00000
pepper_steak 0.52205	wine_market 0.59988	garden_spider 0.58116	dangling_modifier 0.65794
potato_bug 0.50899	cigarette_market 0.59143	mud_dauber 0.56821	cutting_implement 0.63771
tub_thumper 0.49676	tea_market 0.58791	cicada_killer 0.56456	four-wheel-drive_vehicle 0.62591
rum_baba 0.49114	arms_market 0.57372	mining_bee 0.52083	toy_vehicle 0.62213
dice_box 0.48834	sugar_market 0.55928	house_centipede 0.49953	import_system 0.62018
shoo_fly_c -0.02019	tobacco_market_c 0.14430	paper_wasp_c 0.19814	sports_vehicle_c -0.13707



 $p_{i} = W_{i1}u_{1} + W_{i2}u_{2} + W_{i3}u_{3} + W_{i4}v_{1} + W_{i5}v_{2} + W_{i6}v_{3} + b_{i}, \forall i \in \{1,2,3\}$  $p_{i} = W_{i*} \cdot [u; v] + b_{i}, \forall i \in \{1,2,3\}$ 

 $W \in \mathbb{R}^{d \times 2d}$  $u, v, b, p \in \mathbb{R}^{d}$ 



 $p_{i} = W_{i1}u_{1} + W_{i2}u_{2} + W_{i3}u_{3} + W_{i4}v_{1} + W_{i5}v_{2} + W_{i6}v_{3} + b_{i}, \forall i \in \{1,2,3\}$  $p_{i} = W_{i*} \cdot [u; v] + b_{i}, \forall i \in \{1,2,3\}$ 

 $W \in \mathbb{R}^{d \times 2d}$  $u, v, b, p \in \mathbb{R}^{d}$ 



 $p_{i} = W_{i1}u_{1} + W_{i2}u_{2} + W_{i3}u_{3} + W_{i4}v_{1} + W_{i5}v_{2} + W_{i6}v_{3} + b_{i}, \forall i \in \{1,2,3\}$  $p_{i} = W_{i*} \cdot [u; v] + b_{i}, \forall i \in \{1,2,3\}$ 

slobal composition - the two words from any pair of words - e.g. (apple, tree), (car, factory), (hand, puppet) will be composed in the same way

#### **NEIGHBOURS OF MATRIX COMPOSITION**

dog_collar:1	rap_artist:1	property_sale:1	hill_town:1
dog_collar 1.00000	rap_artist 1.00000	property_sale 1.00000	hill_town 1.00000
dog_collar_c 0.54199	rap_artist_c 0.66841	property_sale_c 0.58137	hill_town_c 0.49966
leash 0.48995	rap_star 0.64815	property_transaction 0.54844	mountain_village 0.44740
collar 0.44047	recording_artist 0.56339	property_management 0.52949	market_town 0.40680
dog_tag 0.41794	hip_hop 0.54667	home_sale 0.50604	beach_town 0.37561
shoulder_strap 0.38407	rap 0.54657	estate_agent 0.47905	university_town 0.37451
sports_jacket 0.38068	rap_music 0.52175	property_market 0.46321	market_square 0.37379
brand_manager:3	time_period:3	lake_poets:3	cheese_dip:3
brand_manager 1.00000	time_period $1.00000$	lake_poets 1.00000	cheese_dip 1.00000
marketing_manager 0.70351	time_frame $0.75858$	folk_poet 0.65036	cheese_sauce 0.55767
product_manager 0.65581	period $0.68544$	sand_hopper 0.63909	taco_sauce 0.54454
brand_manager_c 0.64731	time_period_c $0.59635$	lake_poets_c 0.59926	cheese_dip_c 0.51926
marketing_director 0.64690	time_interval $0.59194$	forest_god 0.59435	bechamel_sauce 0.51546
product_marketing 0.62479	year_period $0.58240$	sea_swallow 0.59369	bean_dip 0.50829
marketing_executive 0.59426	time_span $0.57384$	horse_latitude 0.59048	wax_bean 0.50158
rose_hip:1000	$developing\_cost:1000$	snow_job:1000	walkie_talkie:1000
rose_hip 1.00000	developing_cost 1.00000	snow_job 1.00000	walkie_talkie 1.00000
pomegranate 0.54706	morale_building 0.52340	squeeze_play 0.47097	police_radio 0.50134
dog_rose 0.51804	profit_projection 0.51079	snipe_hunt 0.46792	radio_set 0.46869
borage 0.50045	data_formatting 0.49961	tub_thumper 0.44253	cell_phone 0.46803
evening_primrose 0.48132	shareholder_litigation 0.49381	conversation_stopper 0.43753	crystal_set 0.45369
rose_oil 0.47833	inventory_adjustment 0.49292	ham_actor 0.42959	communication_device 0.45242
rose_hip_c 0.18460	developing_cost_c 0.31569	snow_job_c 0.27037	walkie_talkie_c 0.16783

. . . . . . . . . . . . . . . .

. . . . . .

. . . . .

Apfel	Baum	Apfelbaum
'apple'	'tree'	'apple tree'
Holz	Löffel	Holzlöffel
'wood'	'spoon'	'wooden spoon'
Abend	Kleid	Abendkleid
'evening'	'dress'	'evening dress'
Chef	Koch	Chefkoch
'chef'	'cook'	'chef cook'
Mandel	Öl	Mandelöl
'almond'	'oil'	'almond oil'
Orange	Saft	Orangensaft
'orange'	'juice'	'orange juice'
Gast	Bad	Gästebad
'guest'	'bathroom'	'guest bathroom'
Hochzeit	Bild	Hochzeitsbild
'wedding'	'photo'	'wedding photo'

• •

. . . . . . . . . . . .

Apfel	Baum	Apfelbaum
'apple'	'tree'	'apple tree'
Holz	Löffel	Holzlöffel
'wood'	'spoon'	'wooden spoon'
Abend	Kleid	Abendkleid
'evening'	'dress'	'evening dress'
Chef	Koch	Chefkoch
'chef'	'cook'	'chef cook'
Mandel	Öl	Mandelöl
'almond'	'oil'	'almond oil'
Orange	Saft	Orangensaft
'orange'	'juice'	'orange juice'
Gast	Bad	Gästebad
'guest'	'bathroom'	'guest bathroom'
Hochzeit	Bild	Hochzeitsbild
'wedding'	'photo'	'wedding photo'

$u_1$	$v_1$	$\widetilde{p}_1$	$p_1 = f(u_1, v_1)$
$u_2$	$v_2$	$\widetilde{p}_2$	$p_2 = f(u_2, v_2)$
<i>u</i> <sub>3</sub>	<i>v</i> <sub>3</sub>	$\widetilde{p}_3$	$p_3 = f(u_3, v_3)$
$u_4$	$v_4$	$\widetilde{p}_4$	$p_4 = f(u_4, v_4)$
$u_5$	$v_5$	${\widetilde p}_5$	$p_5 = f(u_5, v_5)$
$u_6$	$v_6$	$\widetilde{p}_6$	$p_6 = f(u_6, v_6)$
$u_7$	$v_7$	$\widetilde{p}_7$	$p_7 = f(u_7, v_7)$
<i>u</i> <sub>8</sub>	$v_8$	$\widetilde{p}_8$	$p_8 = f(u_8, v_8)$

Apfel	Baum	Apfelbaum
'apple'	'tree'	'apple tree'
Holz	Löffel	Holzlöffel
'wood'	'spoon'	'wooden spoon'
Abend	Kleid	Abendkleid
'evening'	'dress'	'evening dress'
Chef	Koch	Chefkoch
'chef'	'cook'	'chef cook'
Mandel	Öl	Mandelöl
'almond'	'oil'	'almond oil'
Orange	Saft	Orangensaft
'orange'	'juice'	'orange juice'
Gast	Bad	Gästebad
'guest'	'bathroom'	'guest bathroom'
Hochzeit	Bild	Hochzeitsbild
'wedding'	'photo'	'wedding photo'

$u_1$	$v_1$	$\widetilde{p}_1$	$p_1 = f(u_1, v_1)$
<i>u</i> <sub>2</sub>	$v_2$	$\widetilde{p}_2$	$p_2 = f(u_2, v_2)$
<i>u</i> <sub>3</sub>	<i>v</i> <sub>3</sub>	$\widetilde{p}_3$	$p_3 = f(u_3, v_3)$
$u_4$	$v_4$	$\widetilde{p}_4$	$p_4 = f(u_4, v_4)$
$u_5$	$v_5$	${\widetilde p}_5$	$p_5 = f(u_5, v_5)$
$u_6$	$v_6$	$\widetilde{p}_6$	$p_6 = f(u_6, v_6)$
$u_7$	$v_7$	$\widetilde{p}_7$	$p_7 = f(u_7, v_7)$
<i>u</i> <sub>8</sub>	$v_8$	$\widetilde{p}_8$	$p_8 = f(u_8, v_8)$

Training objective:  $\operatorname{argmax} \frac{p \cdot \tilde{p}}{\|p\|_2 \|\tilde{p}\|_2}$ 

. .

use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

. . . . . . . . . . . . .

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

. . . . . . . . . . . . . .

► split the data into *train-test-dev* portions (70-20-10%)

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- ► split the data into *train-test-dev* portions (70-20-10%)
- ► train the composition model on the *train* data (70%)

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- split the data into train-test-dev portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- ► split the data into *train-test-dev* portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model
    - ► addition:

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- split the data into train-test-dev portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model
    - ► addition: 0 parameters (non-parametric)

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- split the data into train-test-dev portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model
    - ► addition: 0 parameters (non-parametric)
    - ► weighted addition:

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- split the data into train-test-dev portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model
    - ► addition: 0 parameters (non-parametric)
    - ► weighted addition: 2 parameters

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- split the data into train-test-dev portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model
    - ► addition: 0 parameters (non-parametric)
    - ► weighted addition: 2 parameters
    - ► matrix:

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- ► split the data into *train-test-dev* portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model
    - ► addition: 0 parameters (non-parametric)
    - ► weighted addition: 2 parameters
    - ➤ matrix: d x 2d + d parameters (d=200: 200x400+200 = 80,200)

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- ► split the data into *train-test-dev* portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model
    - ► addition: 0 parameters (non-parametric)
    - ► weighted addition: 2 parameters
    - ➤ matrix: d x 2d + d parameters (d=200: 200x400+200 = 80,200)
- choose hyperparameters using the *dev* data (10%) learning rate, optimizer, embedding size, embedding type, etc.

 use as many such triples as you can get (datasets ranging from 4,500 to 240,000 triples)

- ► split the data into *train-test-dev* portions (70-20-10%)
- ► train the composition model on the *train* data (70%)
  - estimate all the parameters of composition model
    - ► addition: 0 parameters (non-parametric)
    - ► weighted addition: 2 parameters
    - ➤ matrix: d x 2d + d parameters (d=200: 200x400+200 = 80,200)
- choose hyperparameters using the *dev* data (10%) learning rate, optimizer, embedding size, embedding type, etc.
- test the composition model with the best hyperparameters <u>once</u> on the *test* data (20%)

### **EVALUATION FOR COMPOSITION MODELS**

intrinsic evaluation - how good are the composed representations created by the model when compared to corpus-derived representations of the same words?

. . . . . . . . . . . . .

• extrinsic evaluation - how good are the composed representations for some external task (e.g. for identifying the semantic relations between words, for parsing, for sentiment analysis etc.)

#### **INTRINSIC EVALUATION**

cosine similarity Apfelbaum 'apple tree' 1 Baum 'tree' 0.97()()()Kirschbaum 'cherry tree'  $\bigcirc$ 0.89  $\bigcirc$ () $\bigcirc$  $\bigcirc$  $\bigcirc$ 0.86 () $\bigcirc$ . . .  $\bigcirc$  $\bigcirc$  $\bigcirc$  $\bigcirc$ . . . . . . Baumstamm 'tree trunk' 0.67  $\bigcirc$  $\bigcirc$  $\bigcirc$ ()Apfel 'apple' 0.65 $\bigcirc$ ()()()() $\bigcirc$  $\bigcirc$ (). . . . . . Schneebesen 'whisk' 0.23  $\bigcirc$  $\bigcirc$  $\bigcirc$ ()Bilderbuch 'picture book' 0.18  $\bigcirc$  $\bigcirc$ () $\bigcirc$ 

#### **INTRINSIC EVALUATION (2)**

\_ \_ \_ Apfelbaum 'apple tree' 1 ()()( ) Baum 'tree'  $\bigcirc$  $\bigcirc$ 0.97  $\bigcirc$ ( ) Kirschbaum 'cherry tree'  $\bigcirc$  $\bigcirc$  $\bigcirc$ 0.89  $\bigcirc$ Apfelbaum 0.88 ← rank 3  $\bigcirc$  $\bigcirc$  $\bigcirc$ 0.86  $\bigcirc$ . . .  $\bigcirc$  $\bigcirc$  $\bigcirc$ ( ) . . . . . . Baumstamm 'tree trunk'  $\bigcirc$  $\bigcirc$  $\bigcirc$ 0.67 ()Apfel 'apple'  $\bigcirc$  $\bigcirc$ 0.65 $\bigcirc$  $\bigcirc$  $\bigcirc$  $\bigcirc$  $\bigcirc$  $\bigcirc$ . . . . . . Schneebesen 'whisk'  $\bigcirc$  $\bigcirc$ 0.23 $\bigcirc$  $\bigcirc$ Bilderbuch 'picture book'  $\bigcirc$ 0.18  $\bigcirc$  $\bigcirc$ ()

cosine similarity

#### **INTRINSIC EVALUATION (3)**

test compound	rank
Telefonkabel 'telephone cable'	1
Schlossturm 'castle tower'	2
Hundekuchen 'dog biscuit'	12
Sitzheizung 'seat heating'	13
Maulwurf 'mole'	1000
Milchmädchen 'milkmaid'	1000



#### **INTRINSIC EVALUATION (3)**

test compound	rank	
Telefonkabel 'telephone cable'	1	
Schlossturm 'castle tower'	2	
		$O_{1}$
Hundekuchen 'dog biscuit'	12	QI
Sitzheizung 'seat heating'	13	are
		of
Maulwurf 'mole'	1000	
Milchmädchen 'milkmaid'	1000	
$1$ $2$ $\cdots$ $12$	13	1000
$\mathbf{Q}_1$	$\mathbf{Q}_2$	$\mathbf{Q}_3$

Q<sub>1</sub>, Q<sub>2</sub>(median), Q<sub>3</sub> are the quartiles of the sorted list of ranks

1000

#### **INTRINSIC EVALUATION (4)**

• •

. . . . . . . . . . . . .

	English Nominal Compounds		
<b>Composition Model</b>	Q1, Q2, Q3	<=5	
Addition	2,7,38	46,14 %	
Weighted Addition	2,7,38	46,14 %	
Matrix	1,2,9	67,37 %	

 $W \in \mathbb{R}^{d \times 2d}$  $u, v, b, p \in \mathbb{R}^{d}$ 



 $p_{i} = W_{i1}u_{1} + W_{i2}u_{2} + W_{i3}u_{3} + W_{i4}v_{1} + W_{i5}v_{2} + W_{i6}v_{3} + b_{i}, \forall i \in \{1,2,3\}$  $p_{i} = W_{i*} \cdot [u; v] + b_{i}, \forall i \in \{1,2,3\}$ 

slobal composition - the two words from any pair of words - e.g. (apple, tree), (car, factory), (hand, puppet) will be composed in the same way

#### FULLEX COMPOSITION – SOCHER ET AL., 2012 (1)



#### FULLEX COMPOSITION – SOCHER ET AL., 2012 (1)

 $u_V, v_U, b, p \in \mathbb{R}^d$ 

 $W \in \mathbb{R}^{d \times 2d}$ 



#### FULLEX COMPOSITION – SOCHER ET AL., 2012 (2)



#### FULLEX COMPOSITION – SOCHER ET AL., 2012 (2)

u = vector for apple



#### FULLEX COMPOSITION – SOCHER ET AL., 2012 (2)

u = vector for apple

v = vector for tree


### FULLEX COMPOSITION – SOCHER ET AL., 2012 (2)

*u* = *vector for* apple

v = vector for tree



 $u, u_V, v, v_U \in \mathbb{R}^d$  $V, U \in \mathbb{R}^{d \times d}$  $T \in \mathbb{R}^{|V| \times d \times d}$ 

## FULLEX COMPOSITION – SOCHER ET AL., 2012 (2)

u = vector for apple

v = vector for tree



### FULLEX COMPOSITION – SOCHER ET AL., 2012 (3)



 $p_i = W_{i*} \cdot [u_V; v_U] + b_i, \forall i \in \{1, 2, 3\}$ 

### **NEIGHBOURS OF FULLLEX COMPOSITION**

monitor_program:1	party_chief:1	rogue_elephant:1	service_division:1
monitor_program 1.00000	party_chief 1.00000	rogue_elephant 1.00000	service_division 1.00000
monitor_program_c 0.48685	party_chief_c 0.57671	rogue_elephant_c 0.66267	service_division_c 0.70384
import_system 0.44099	party_boss 0.45157	horse_doctor 0.66240	service_department 0.56786
program_quality 0.43015	party_official 0.43346	plow_horse 0.65587	service_facility 0.56749
traffic_mitigation 0.42661	party_leader 0.42197	chicken_yard 0.65222	service_operation 0.51915
personality_quiz 0.42372	politburo_member 0.41440	bank_guard 0.63444	service_program 0.51402
buffer_store 0.42320	party_branch 0.41280	government_man 0.62989	service_office 0.51309
mailsystem:2	vegetable_garden:2	$shooting_lodge:2$	guide_rope:2
mail_system 1.00000	vegetable_garden 1.00000	shooting_lodge $1.00000$	guide_rope 1.00000
email_system 0.64818	vegetable_patch 0.76619	hunting_lodge $0.61701$	center_pole 0.68092
mail_system_c 0.60377	vegetable_garden_c 0.74822	shooting_lodge_c $0.51058$	guide_rope_c 0.64363
voice_mail 0.55972	kitchen_garden 0.73704	shooting_box $0.50183$	grappling_iron 0.64078
mail_service 0.54594	herb_garden 0.71074	cow_barn $0.47860$	water_wagon 0.61621
mail_box 0.45023	flower_bed 0.67242	army_hut $0.46117$	snatch_block 0.60451
snail_mail 0.42475	fruit_tree 0.64774	manor_house $0.45763$	starting_stall 0.60083
custom_officer:1000	sounding_board:1000	service_break:1000	home_plate:1000
custom_officer 1.00000	sounding_board 1.00000	service_break 1.00000	home_plate 1.00000
custom_official 0.76468	rallying_point 0.42654	rating_period 0.60125	plate_umpire 0.68733
stretcher_party 0.71902	launching_pad 0.41923	chow_line 0.56250	center_field 0.49705
slip_coach 0.69349	testing_ground 0.41371	operation_officer 0.55321	umpire 0.49245
riot_officer 0.68981	reference_point 0.38087	pace_lap 0.55097	dugout 0.46930
poll_taker 0.68704	teaching_tool 0.37899	moving_staircase 0.54872	strike_zone 0.45223
custom_officer_c 0.43134	sounding_board_c 0.15510	service_break_c 0.35390	home_plate_c 0.09510

► the number of parameters increases with the size of the vocabulary

 $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200

- $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
  - ► T size: 100 x 200 x 200 = 100 x 40,000 = 4,000,000

- $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
  - ► T size: 100 x 200 x 200 = 100 x 40,000 = 4,000,000
  - ► W size: 200 x 400 + 200 = 80,200

- $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
  - ► T size: 100 x 200 x 200 = 100 x 40,000 = 4,000,000
  - ► W size: 200 x 400 + 200 = 80,200
  - ► total: 4,080,200 parameters

- $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
  - ► T size: 100 x 200 x 200 = 100 x 40,000 = 4,000,000
  - ► W size: 200 x 400 + 200 = 80,200
  - ► total: 4,080,200 parameters
- ► |V| = 10,000 words; d=200

- $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
  - ► T size: 100 x 200 x 200 = 100 x 40,000 = 4,000,000
  - ► W size: 200 x 400 + 200 = 80,200
  - ► total: 4,080,200 parameters
- ► |V| = 10,000 words; d=200
  - ► T size:

- ► |V| = 100 words; d (size of word representation) = 200
  - ► T size: 100 x 200 x 200 = 100 x 40,000 = 4,000,000
  - ► W size: 200 x 400 + 200 = 80,200
  - ► total: 4,080,200 parameters
- ► |V| = 10,000 words; d=200
  - ► T size:
    - ► 10,000 x 200 x 200 = 10,000 x 40,000 = 400,000,000

- ► |V| = 100 words; d (size of word representation) = 200
  - ► T size: 100 x 200 x 200 = 100 x 40,000 = 4,000,000
  - ► W size: 200 x 400 + 200 = 80,200
  - ► total: 4,080,200 parameters
- ► |V| = 10,000 words; d=200
  - ► T size:
    - ► 10,000 x 200 x 200 = 10,000 x 40,000 = 400,000,000
  - ► W size: 200 x 400 + 200 = 80,200

- $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
  - ► T size: 100 x 200 x 200 = 100 x 40,000 = 4,000,000
  - ► W size: 200 x 400 + 200 = 80,200
  - ► total: 4,080,200 parameters
- ► |V| = 10,000 words; d=200
  - ► T size:
    - ► 10,000 x 200 x 200 = 10,000 x 40,000 = 400,000,000
  - ► W size: 200 x 400 + 200 = 80,200
  - ► total: 400,080,000 parameters



Ι.

 $u, v, u^m, v^h, u \odot u^m, v \odot v^h \in \mathbb{R}^d$ 



Ι.

 $u, v, u^m, v^h, u \odot u^m, v \odot v^h \in \mathbb{R}^d$  $u_1 u_1^m$  $v_1 v_1^h$  $u_1$  $v^h$  $u \odot u^m$  $\begin{bmatrix} u_2^m \\ u_3^m \end{bmatrix}$  $v_2^h$  $v_3^h$  $v \odot v^h$ <u>u</u> 📀 v<sub>2</sub>  $u_2 u_2^m$  $v \odot$  $v_2 v_2^h$ *u*<sub>2</sub>  $v_3$  $u_3 u_3^m$  $v_3 v_3^h$  $u_3$  $a_1^m$  $u_1^m$  $a_1^h$  $u_1^h$  $v_1^m$  $v_1^h$  $W_M \in \mathbb{R}^{d \times |V|}$  $W_H \in \mathbb{R}^{d \times |V|}$  $u_2^m$  $a_2^m$  $|v_2^m|$  $a_2^h$  $u_2^h$  $v_2^h$ . . . • • • . . . . . .  $v_3^m$  $a_3^m$  $u_3^m$  $a_3^h$  $u_3^h$  $v_3^h$ factory factory company company car car vector masks are initialized

with a vector of ones

 $u, v, u^m, v^h, u \odot u^m, v \odot v^h \in \mathbb{R}^d$  $v_2^h$  $v_3^h$  $-u \odot u^m$  $u_2^m$  $u_3^m$  $v \odot v^h$ <u>u</u> 📀 **v**<sub>2</sub>  $\succ v \odot$  $u_2 u_2^m$  $v_2 v_2^h$  $u_2$  $v_3$  $u_{3}u_{3}^{m}$  $v_3 v_3^h$  $u_3$  $u_1^m$  $u_1^h$  $W_M \in \mathbb{R}^{d \times |V|}$  $|u_2^m|$  $W_H \in \mathbb{R}^{d \times |V|}$  $a_2^m$  $v_2^m$  $a_2^h$ •••  $\left[u_{2}^{h}\right]$  $v_2^h$  $a_3^m$  $v_3^m$  $u_3^m$  $a_3^h$  $u_3^h$  $v_3^h$ factory factory company company car car vector masks are initialized with a vector of ones

*car factory*:

Ι.

 $u = vector for car; u^m = vector mask for modifier car$  $v = vector for factory; v^h = vector mask for head factory$ 

 $u, v, u^m, v^h, u \odot u^m, v \odot v^h \in \mathbb{R}^d$  $> u \odot u^m$  $v_2^h$  $v_3^h$  $u_2^m$  $u_3^m$ **v**<sub>2</sub>  $> v \odot$ <u>u</u> 💽  $v \odot v^h$  $u_2 u_2^m$  $v_2 v_2^h$ = $u_2$  $u_3 u_3^m$  $v_3 v_3^h$  $u_3$  $W_M \in \mathbb{R}^{d \times |V|}$  $|u_2^m|$  $W_H \in \mathbb{R}^{d \times |V|}$ **a**<sup>*h*</sup>  $|a_2^m|$ •••  $\left[u_{2}^{h}\right]$  $a_3^m$  $u_3^h$ factory company factory company car car vector masks are initialized with a vector of ones

car factory:

Ι.

 $u = vector for car; u^m = vector mask for modifier car$  $v = vector for factory; v^h = vector mask for head factory$ 

company car:

 $u = vector for ?; u^m = vector mask for modifier ?$  $v = vector for ?; v^h = vector mask for head ?$ 

П.







 $WMask: p_i = W_{i*} \cdot [u \odot u^m; v \odot v^h] + b_i, \forall i \in \{1,2,3\}$  $AddMask: p_i = u \odot u^m + v \odot v^h$ 

### **NEIGHBOURS OF WMASK COMPOSITION**

disk_error:1	currency_value:1	employee_manual:1	metal_armor:1
disk_error $1.00000$	currency_value 1.00000	employee_manual 1.00000	metal_armor 1.00000
disk_error_c $0.63742$	currency_value_c 0.68259	employee_manual_c 0.69876	metal_armor_c 0.68537
hardware_error $0.62881$	currency_rate 0.61195	company_attorney 0.59591	chain_armour 0.67611
program_error $0.56904$	exchange_rate 0.54930	apple_fritter 0.58195	face_guard 0.65524
system_error $0.54074$	currency_depreciation 0.53445	redundancy_policy 0.57816	bicycle_clip 0.64259
head_crash $0.52071$	equity_price 0.49077	intelligence_cell 0.57682	bathing_trunks 0.63836
computer_error $0.51923$	currency_policy 0.49042	x-ray_scan 0.57464	bathing_cap 0.61547
church_leader:3	majority_party:3	picnic_table:3	sweetheart_deal:3
church_leader 1.00000	majority_party 1.00000	picnic_table 1.00000	sweetheart_deal 1.00000
christian_leader 0.74565	minority_party 0.73798	fire_pit 0.72293	backroom_deal 0.65540
church_member 0.69506	party_caucus 0.60063	picnic_area 0.69909	barter_deal 0.54495
church_leader_c 0.64688	majority_party_c 0.49051	<b>picnic_table_c 0.60321</b>	sweetheart_deal_c 0.54421
church_official 0.61824	opposition_party 0.47730	picnic 0.58632	oil_deal 0.51831
clergy 0.61303	party_leader 0.46938	barbecue_pit 0.55666	business_deal 0.51084
community_leader 0.56897	majority_coalition 0.46083	basketball_court 0.54764	marketing_deal 0.48358
screen_price:1000	police_lineup:1000	fossil_oil:1000	deed_poll:1000
screen_price 1.00000	police_lineup 1.00000	fossil_oil 1.00000	deed_poll 1.00000
coach_fare 0.56077	bank_guard 0.57622	import_system 0.68093	marriage_certificate 0.53045
one-off_charge 0.53977	mopping-up_operation 0.57314	banana_export 0.65954	family_name 0.46612
manufacturing_level 0.51958	police_sweep 0.57112	price_distribution 0.64716	company_name 0.44564
seat_price 0.51609	surprise_inspection 0.56619	self_loader 0.64239	birth_certificate 0.40850
carriage_charge 0.51505	trash_barrel 0.56103	trading_transaction 0.64030	given_name 0.39348
screen_price_c 0.31463	police_lineup_c 0.32903	fossil_oil_c 0.37813	deed_poll_c 0.10163

. . . . . .

the number of parameters still increases with the size of the vocabulary

. . . . . . .

. . . . . . . .

. . . . . . . . .

- the number of parameters still increases with the size of the vocabulary
  - ► |V| = 100 words; d (size of word representation) = 200

- the number of parameters still increases with the size of the vocabulary
  - $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
    - ►  $W_M$ ,  $W_H$  size: 100 x (200 + 200) = 100 x 400 = 40,000

- the number of parameters still increases with the size of the vocabulary
  - $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
    - ► W<sub>M</sub>, W<sub>H</sub> size: 100 x (200 + 200) = 100 x 400 = 40,000
    - ► W size: 200 x 400 + 200 = 80,200

- the number of parameters still increases with the size of the vocabulary
  - $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
    - ► W<sub>M</sub>, W<sub>H</sub> size: 100 x (200 + 200) = 100 x 400 = 40,000
    - ► W size: 200 x 400 + 200 = 80,200
    - ► total: 120,200 parameters

- the number of parameters still increases with the size of the vocabulary
  - $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
    - ► W<sub>M</sub>, W<sub>H</sub> size: 100 x (200 + 200) = 100 x 400 = 40,000
    - ► W size: 200 x 400 + 200 = 80,200
    - ► total: 120,200 parameters
  - ► |V| = 10,000 words; d=200

- the number of parameters still increases with the size of the vocabulary
  - $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
    - ► W<sub>M</sub>, W<sub>H</sub> size: 100 x (200 + 200) = 100 x 400 = 40,000
    - ► W size: 200 x 400 + 200 = 80,200
    - ► total: 120,200 parameters
  - ► |V| = 10,000 words; d=200
    - ► W<sub>M</sub>, W<sub>H</sub> size:

- the number of parameters still increases with the size of the vocabulary
  - $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
    - ► W<sub>M</sub>, W<sub>H</sub> size: 100 x (200 + 200) = 100 x 400 = 40,000
    - ► W size: 200 x 400 + 200 = 80,200
    - ► total: 120,200 parameters
  - ► |V| = 10,000 words; d=200
    - ► W<sub>M</sub>, W<sub>H</sub> size:
      - ► 10,000 x (200 + 200) = 10,000 x 400 = 4,000,000

- the number of parameters still increases with the size of the vocabulary
  - $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
    - ► W<sub>M</sub>, W<sub>H</sub> size: 100 x (200 + 200) = 100 x 400 = 40,000
    - ► W size: 200 x 400 + 200 = 80,200
    - ► total: 120,200 parameters
  - ► |V| = 10,000 words; d=200
    - ► W<sub>M</sub>, W<sub>H</sub> size:
      - ► 10,000 x (200 + 200) = 10,000 x 400 = 4,000,000
    - ► W size: 200 x 400 + 200 = 80,200
# MASK MODELS ALLEVIATE THE SIZE ISSUE OF FULLLEX

- the number of parameters still increases with the size of the vocabulary
  - $\blacktriangleright$  |V| = 100 words; d (size of word representation) = 200
    - ► W<sub>M</sub>, W<sub>H</sub> size: 100 x (200 + 200) = 100 x 400 = 40,000
    - ► W size: 200 x 400 + 200 = 80,200
    - ► total: 120,200 parameters
  - ► |V| = 10,000 words; d=200
    - ► W<sub>M</sub>, W<sub>H</sub> size:
      - ► 10,000 x (200 + 200) = 10,000 x 400 = 4,000,000
    - ► W size: 200 x 400 + 200 = 80,200
    - ► total: 4,080,000 parameters

### **INTRINSIC EVALUATION (5)**

. . . . . . . .

. . .

. . . . . . . . . . . . .

	English Nominal Compounds	
<b>Composition Model</b>	Q1, Q2, Q3	<=5
Addition	2,7,38	46,14 %
Weighted Addition	2,7,38	46,14 %
Matrix	1,2,9	67,37 %
WMask	1,2,7	71,53 %
FullLex	1,2,7	72,82 %

the word matrices/masks are "trained" only for those words in the training data

- the word matrices/masks are "trained" only for those words in the training data
- for words that are not seen during training, fulllex/mask models are equivalent to

- the word matrices/masks are "trained" only for those words in the training data
- for words that are not seen during training, fulllex/mask models are equivalent to
  - ► the matrix model!

. . . . . . . . . . .

the word matrices/masks are individual - no sharing of parameters between similar words: e.g. 'blue dress' 'sky-blue dress'

- the word matrices/masks are individual no sharing of parameters between similar words: e.g. 'blue dress' 'sky-blue dress'
- if words are rare in training, their matrices/masks will get less updates

- the word matrices/masks are individual no sharing of parameters between similar words: e.g. 'blue dress' 'sky-blue dress'
- if words are rare in training, their matrices/masks will get less updates
- they will remain closer to the identity matrices/one vector, and therefore have less of an impact on the final representation

the word matrices/masks are individual - no sharing of parameters between similar words: e.g. 'blue dress' 'sky-blue dress'

- if words are rare in training, their matrices/masks will get less updates
- they will remain closer to the identity matrices/one vector, and therefore have less of an impact on the final representation
- the amount of updates that the matrix/mask of each word will get depends on its train frequency

the word matrices/masks are individual - no sharing of parameters between similar words: e.g. 'blue dress' 'sky-blue dress'

- if words are rare in training, their matrices/masks will get less updates
- they will remain closer to the identity matrices/one vector, and therefore have less of an impact on the final representation
- the amount of updates that the matrix/mask of each word will get depends on its train frequency
- cannot use the fact that 'blue' and 'sky-blue' are colors

#### REFERENCES

- Jeff Mitchell and Mirella Lapata. 2010. Composition in Distributional Models of Semantics. Cognitive Science 34 (2010) 1388–1429.
- Richard Socher, Christopher Manning and Andrew Ng. 2010. Learning Continuous Phrase Representations and Syntactic Parsing with Recursive Neural Networks. In Proceedings of NIPS-2010 Deep Learning Workshop.
- Richard Socher, Brody Huval, Christopher Manning and Andrew Ng. 2012. Semantic Compositionality through Recursive Matrix-Vector Spaces. In Proceedings of EMNLP-CoNLL 2012.
- Corina Dima. 2015. Reverse-engineering Language: A Study on the Semantic Compositionality of German Compounds. In Proceedings of EMNLP 2015.

### **IMAGE CREDITS**

<u>https://pixabay.com/photos/code-geek-talk-code-to-me-coffee-</u> <u>cup-2680204/</u>

https://commons.wikimedia.org/wiki/File:Hand.svg

https://en.wikipedia.org/wiki/File:Kewpie\_doll.jpg

https://bar.wikipedia.org/wiki/Datei:Hand\_Puppet.jpg

https://de.wikipedia.org/wiki/Datei:Theater\_Baden-Baden-22-gje.jpg

<u>https://commons.wikimedia.org/wiki/</u> <u>File:Bamboccianti c1700 Theaterauff%C3%BChrung im Freien.jpg</u>

<u>https://commons.wikimedia.org/wiki/</u> <u>File:Bundesarchiv\_Bild\_183-76973-0003, Leipzig, Opernhaus, Auff%</u> <u>C3%BChrung\_%22Dornr%C3%B6schen%22.jpg</u>

https://pixabay.com/de/photos/kasperle-puppentheater-puppekasper-677526/