

# A word is worth a **thousand vectors**

(word2vec, lda, and introducing lda2vec)

Christopher Moody  
@ Stitch Fix

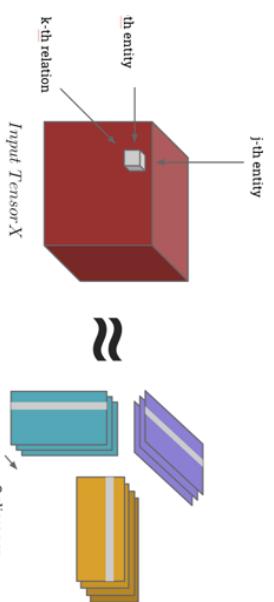
# About

Gaussian Processes

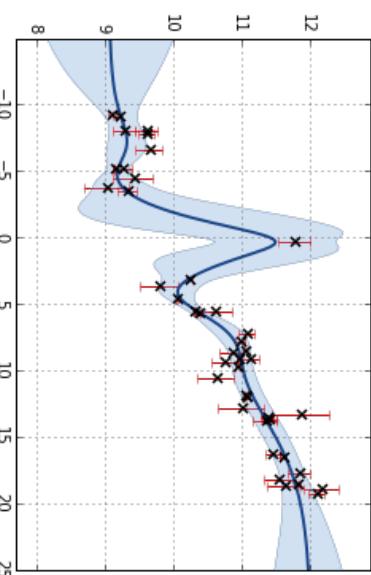
t-SNE



## Tensor Decomposition



chainer  
deep learning



@chrismoody  
Caltech Physics

PhD. in astrostats supercomputing

Sklearn t-SNE contributor

Data Labs at Stitch Fix

[github.com/cemoody](https://github.com/cemoody)



# Credit

Large swathes of this talk are from previous presentations by:

- Tomas Mikolov
- David Blei
- Christopher Olah
- Radim Rehurek
- Omer Levy & Yoav Goldberg
- Richard Socher
- Xin Rong
- Tim Hopper

1

word2vec

2  
lda

3  
lda2vec

## word2vec

1.  $king - man + woman = queen$
2. Huge splash in NLP world
3. Learns from raw text
4. Pretty simple algorithm
5. Comes pretrained

# word2vec

1. Set up an objective function
2. Randomly initialize vectors
3. Do gradient descent

word2vec: learn word vector  $v_{in}$   
from it's surrounding context

$v_{in}$

“The fox jumped over the lazy dog”

Maximize the likelihood of seeing the words given the word **over**.

$$P(\text{the}|\text{over})$$

$$P(\text{fox}|\text{over})$$

$$P(\text{jum}ped|\text{over})$$

$$P(\text{the}|\text{over})$$

$$P(\text{laz}y|\text{over})$$

$$P(\text{dog}|\text{over})$$

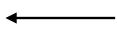
...instead of maximizing the likelihood of co-occurrence counts.

What should this be?

$$P(\text{fox}|\text{over})$$

Should depend on the word vectors.

$$P(\text{fox}|\text{over})$$



$$P(v_{\text{fox}}|v_{\text{over}})$$

Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”

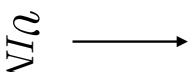
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



$v_{IN}$

Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



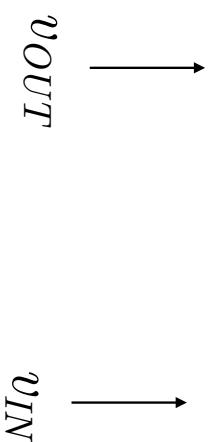
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



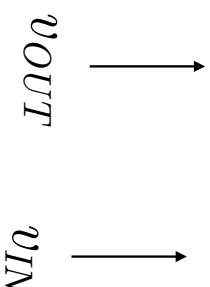
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



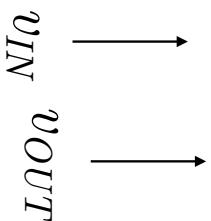
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



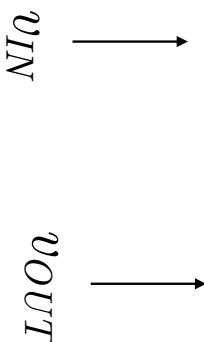
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



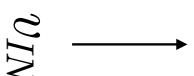
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over **the** lazy dog”



Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



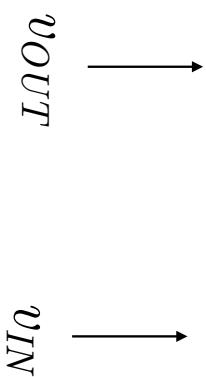
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



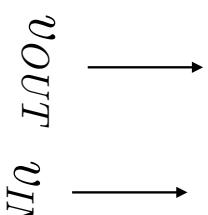
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



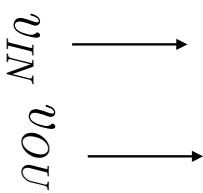
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



$v_{IN}$   $v_{OUT}$

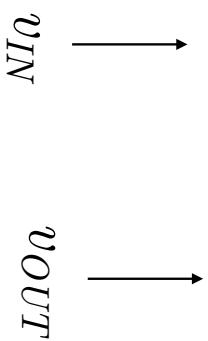
Twist: we have *two* vectors for every word.

Should depend on whether it's the input or the output.

Also a *context* window around every input word.

$$P(v_{OUT}|v_{IN})$$

“The fox jumped over the lazy dog”



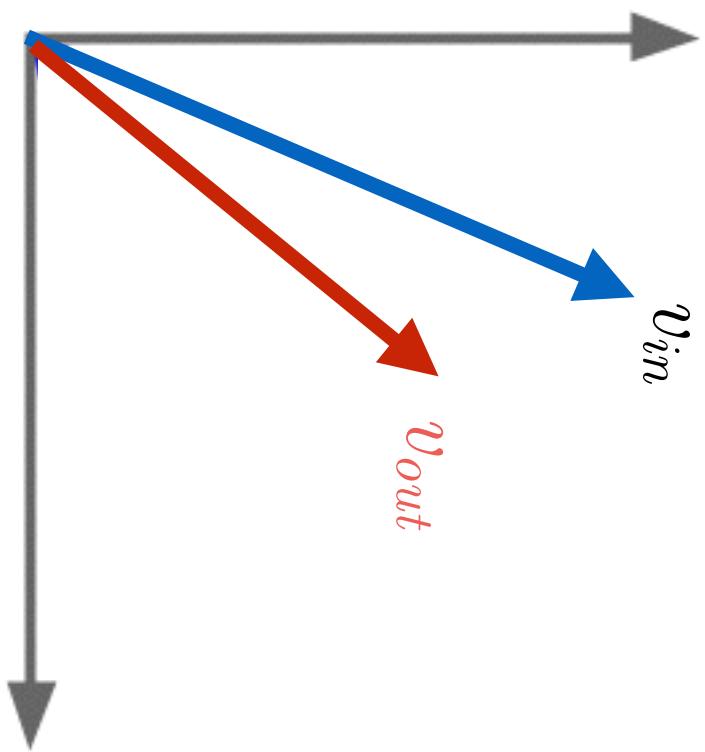
How should we define  $P(v_{OUT}|v_{IN})$ ?

Measure loss between

$v_{IN}$  and  $v_{OUT}$ ?

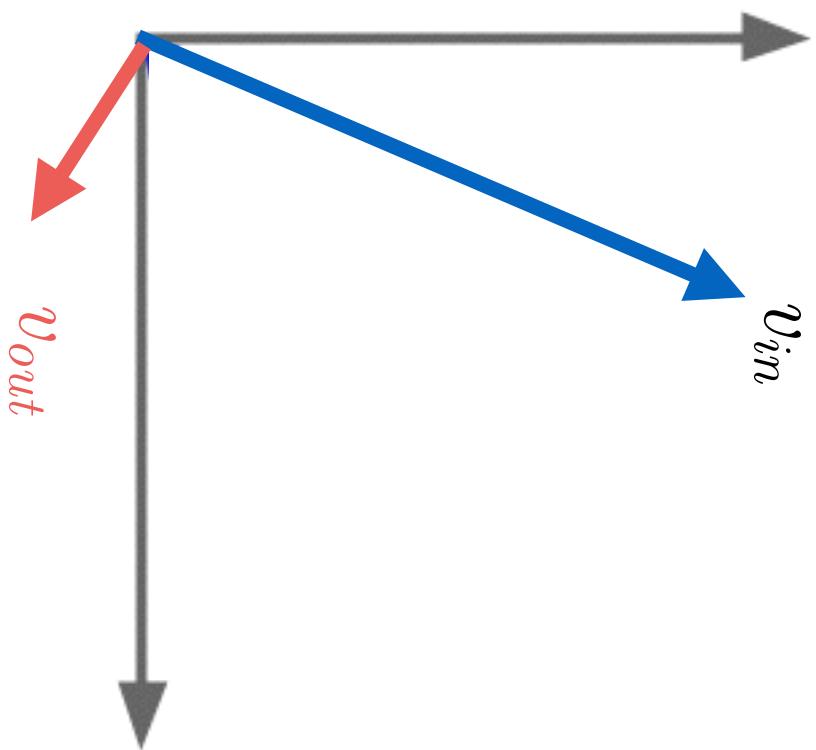
$v_{in} \bullet v_{out}$

objective



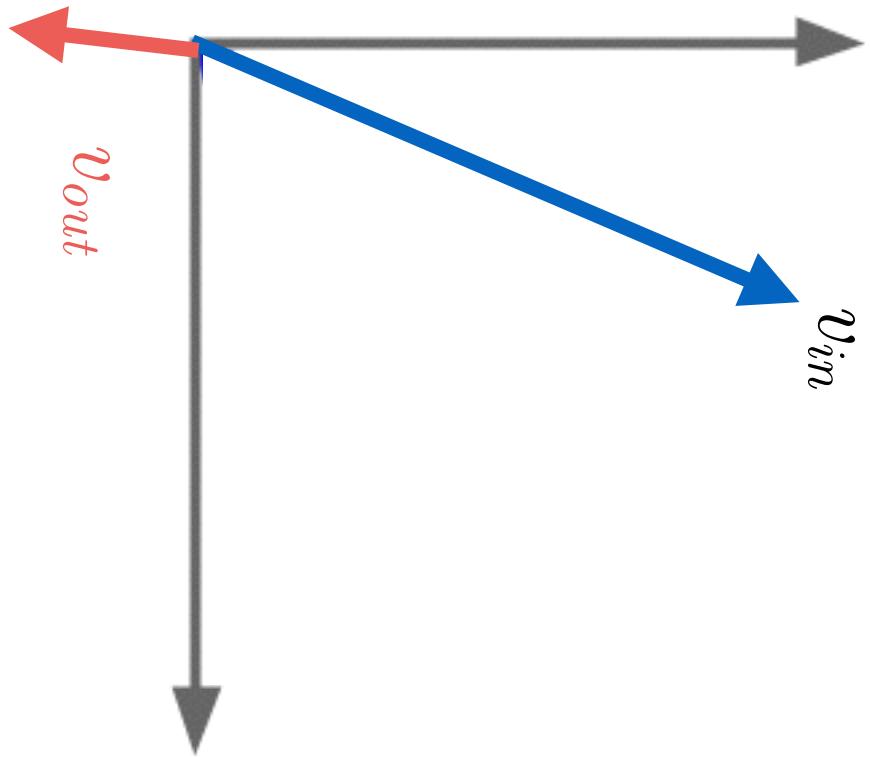
$$v_{in} \cdot v_{out} \sim 1$$

objective



$$v_{in} \cdot v_{out} \sim 0$$

objective



$$v_{in} \cdot v_{out} \approx -1$$

objective

$$v_{in} \bullet v_{out} \in [-1,1]$$

objective

But we'd like to measure a probability.

$$v_{in} \bullet v_{out} \in [-1,1]$$

But we'd like to measure a probability.

$$\text{softmax}(v_{in} \bullet v_{out}) \in [0,1]$$

But we'd like to measure a probability.

$$\text{softmax}(v_{in} \bullet v_{out})$$

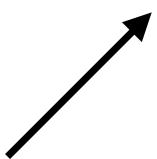
Probability of choosing 1 of N discrete items.  
Mapping from vector space to a multinomial over words.

But we'd like to measure a probability.

$$\text{softmax} \sim \exp(v_{in} \bullet v_{out}) \in [0,1]$$

But we'd like to measure a probability.

$$\text{softmax} = \frac{\exp(v_{in} \bullet v_{out})}{\sum_{k \in V} \exp(v_{in} \bullet v_k)}$$



Normalization term over all words

But we'd like to measure a probability.

$$\text{softmax} = \frac{\exp(v_{in} \bullet v_{out})}{\sum_{k \in V} \exp(v_{in} \bullet v_k)} = P(v_{out} | v_{in})$$

Learn by gradient descent on the softmax prob.

For every example we see update  $v_{in}$

$$v_{in} := v_{in} + \frac{\partial}{\partial v_{in}} P(v_{out} | v_{in})$$

$$v_{out} := v_{out} + \frac{\partial}{\partial v_{out}} P(v_{out} | v_{in})$$

Model (training time)	Redmond	Havel	ninjutsu
Collobert (50d) (2 months)	conyers lubbock keene	plauen dzerzhinsky osterreich	reiki kohana karate
Turian (200d) (few weeks)	McCarthy Alston Cousins	Jewell Arzu Ovitz	- - -
Mnih (100d) (7 days)	Podhurst Harlang Agarwal	Pontiff Pinochet Rodionov	- - -
Skip-Phrase (1000d, 1 day)	Redmond Wash. Redmond Washington Microsoft	Vaclav Havel president Vaclav Havel Velvet Revolution	ninja martial arts swordsmanship

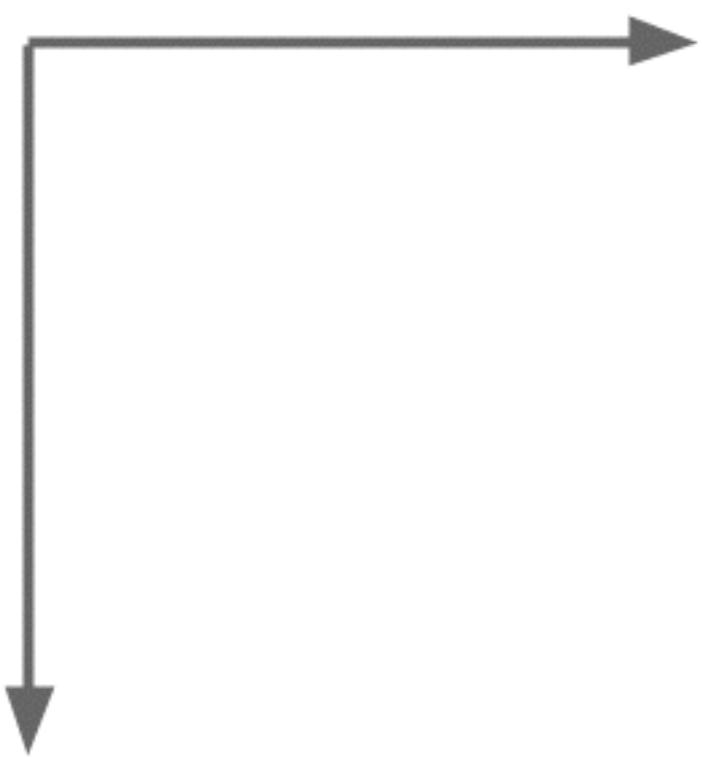


Word2vec

Model	Vector Dimensionality	Training words	Accuracy [%]		
			Semantic	Syntactic	Total
Collobert-Weston NNLM	50	660M	9.3	12.3	11.0
Turian NNLM	50	37M	1.4	2.6	2.1
Turian RNLM	200	37M	1.4	2.2	1.8
Mnih NNLM	50	37M	1.8	9.1	5.8
Mnih NNLM	100	37M	3.3	13.2	8.8
Mikolov RNNLM	80	320M	4.9	18.4	12.7
Mikolov RNNLM	640	320M	8.6	36.5	24.6
Huang NNLM	50	990M	13.3	11.6	12.3
CBOW	300	783M	15.5	53.1	36.1
Skip-gram	300	783M	<b>50.0</b>	55.9	<b>53.3</b>

← WORD2VEC

What is king + man - woman?



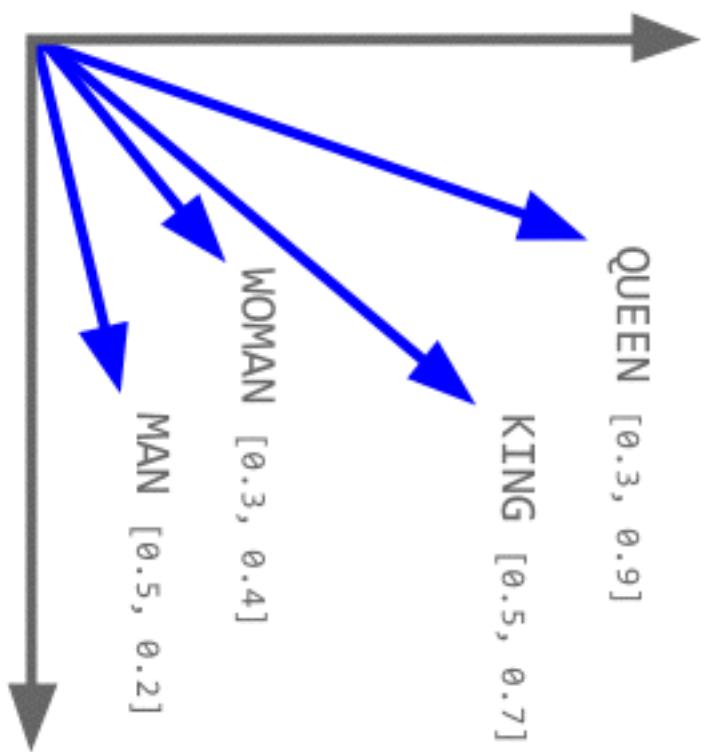
Load up the word vectors

QUEEN [0.3, 0.9]

KING [0.5, 0.7]

WOMAN [0.3, 0.4]

MAN [0.5, 0.2]

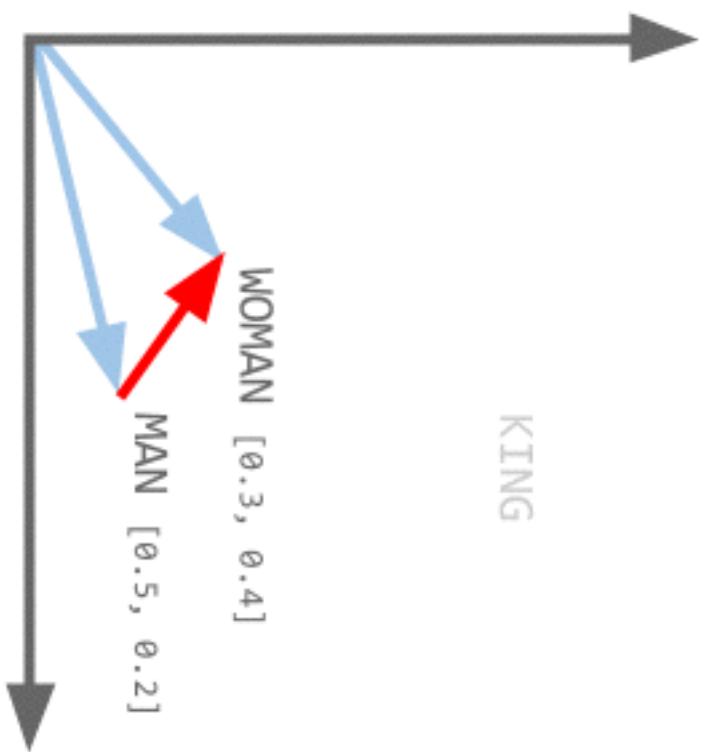


Start with man – woman

KING

WOMAN [0.3, 0.4]

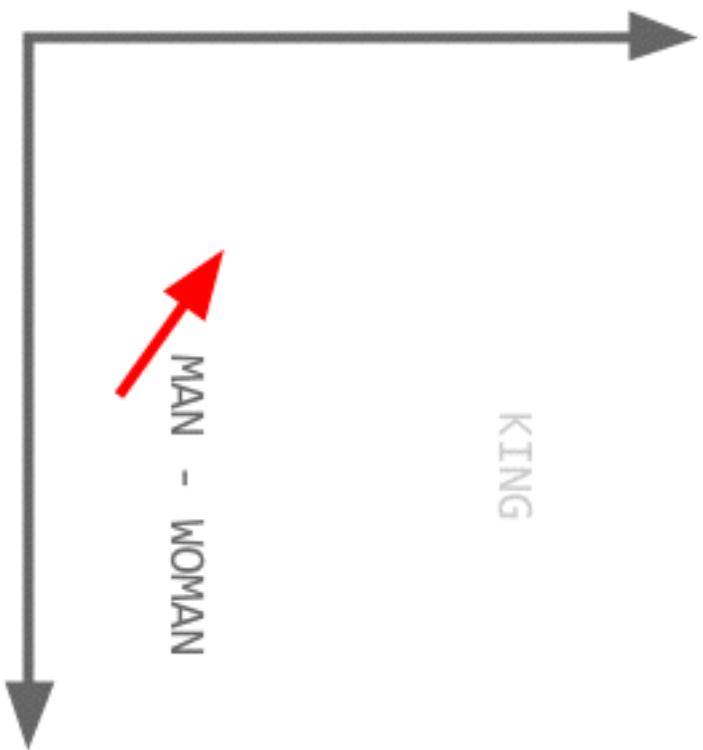
MAN [0.5, 0.2]



Start with man - woman

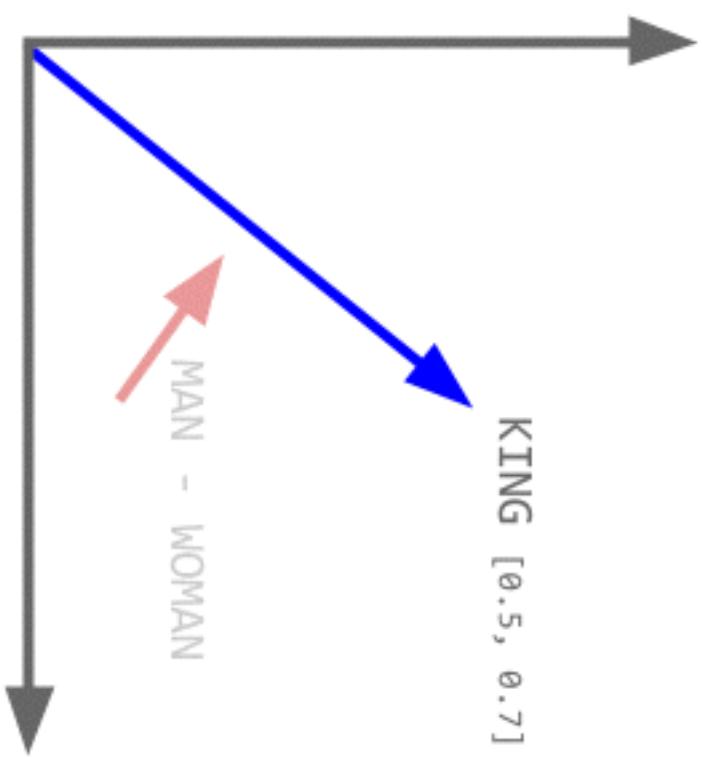
KING

MAN - WOMAN

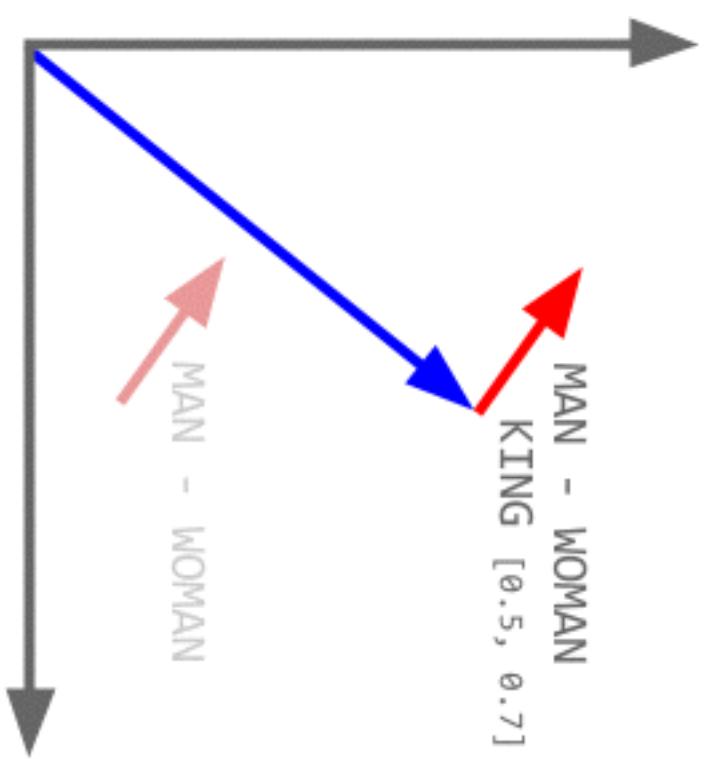


Then take king

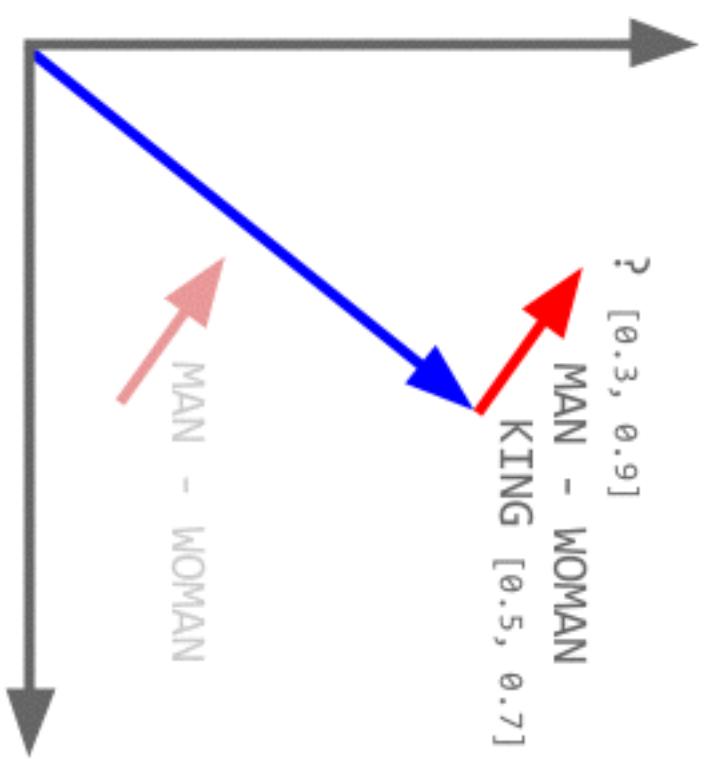
KING [0.5, 0.7]



And add man - woman



And add man - woman

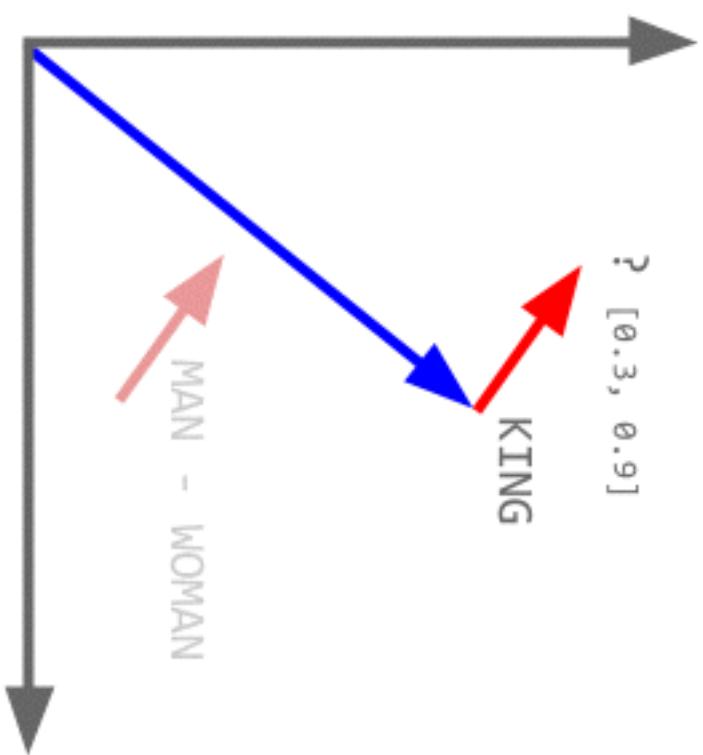


Find nearest word to result

? [0.3, 0.9]

KING

MAN - WOMAN

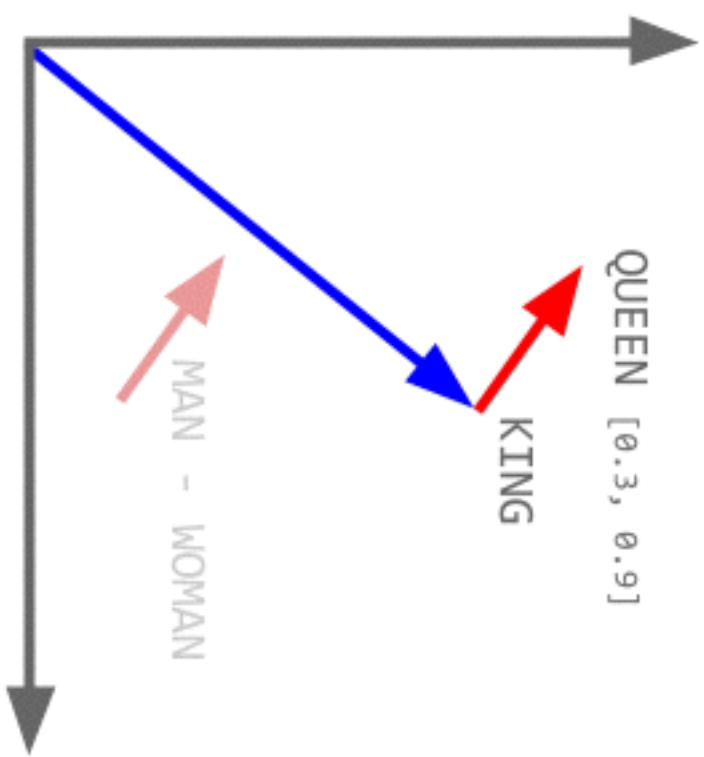


queen is closest to resulting vector

QUEEN [0.3, 0.9]

KING

MAN - WOMAN

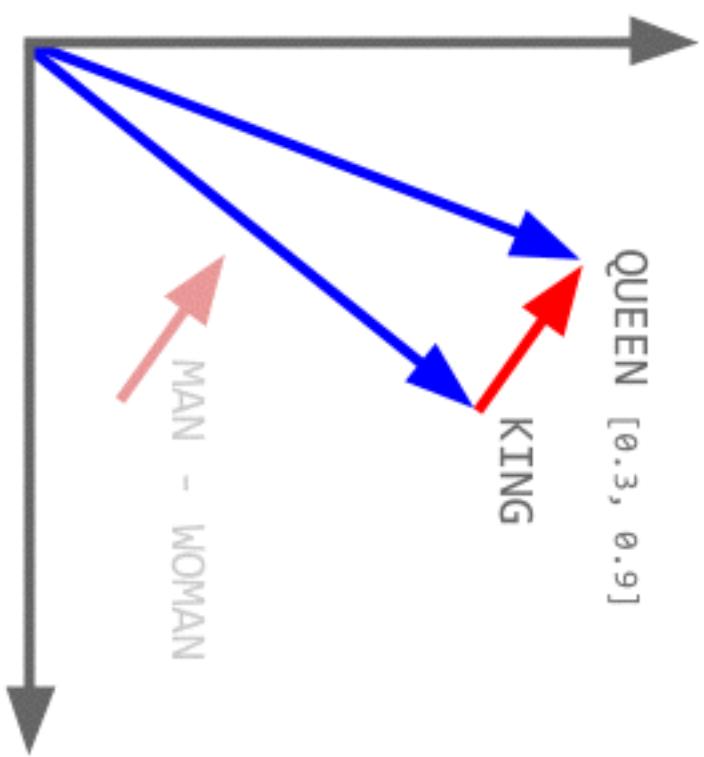


queen is closest to resulting vector

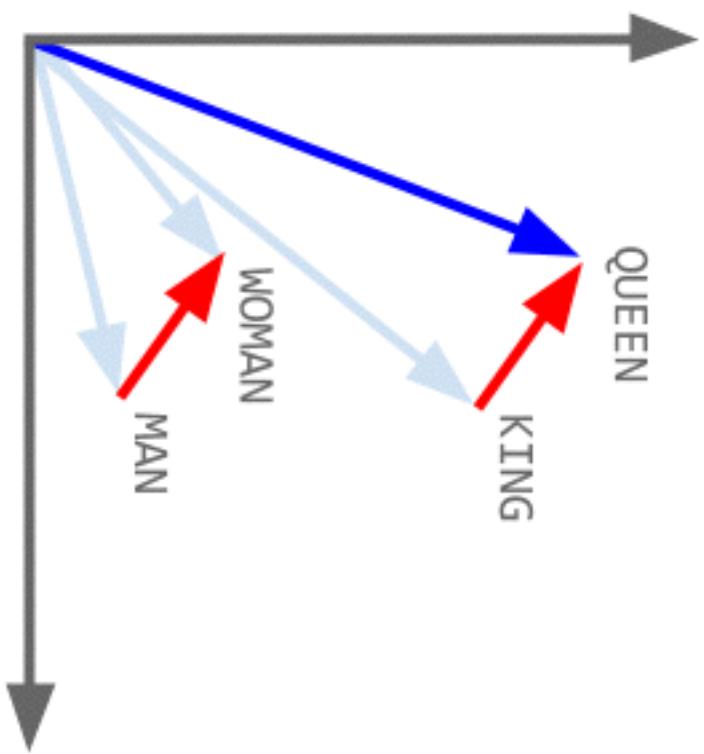
QUEEN [0.3, 0.9]

KING

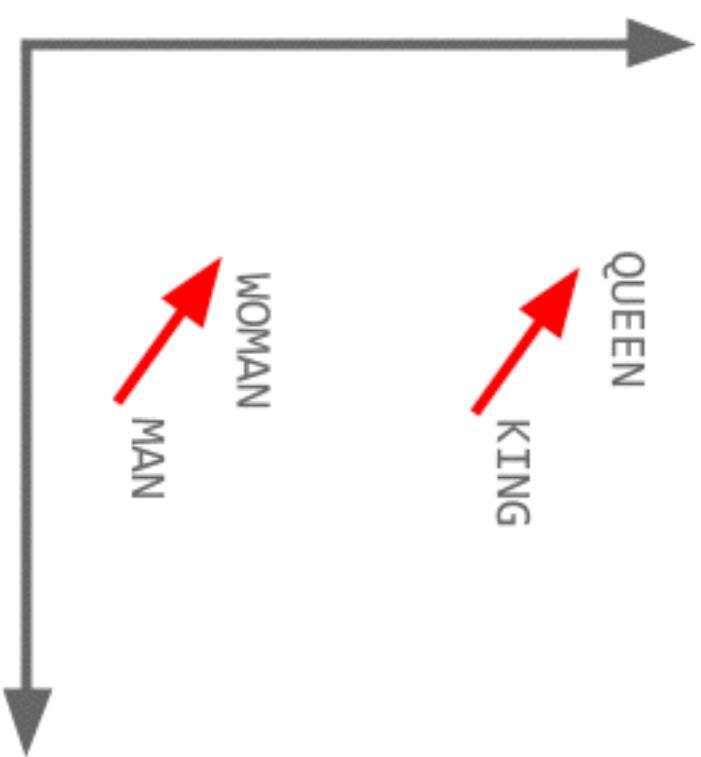
MAN - WOMAN



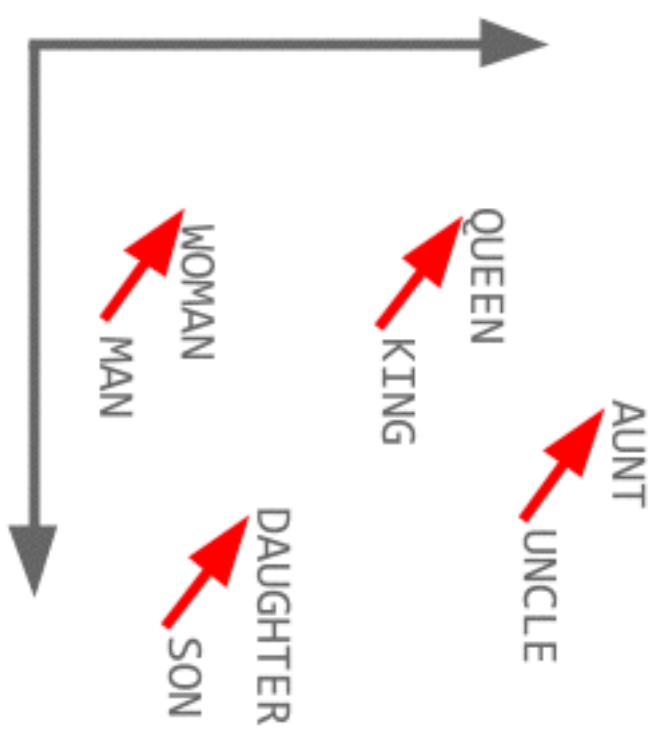
So king + man - woman = queen!



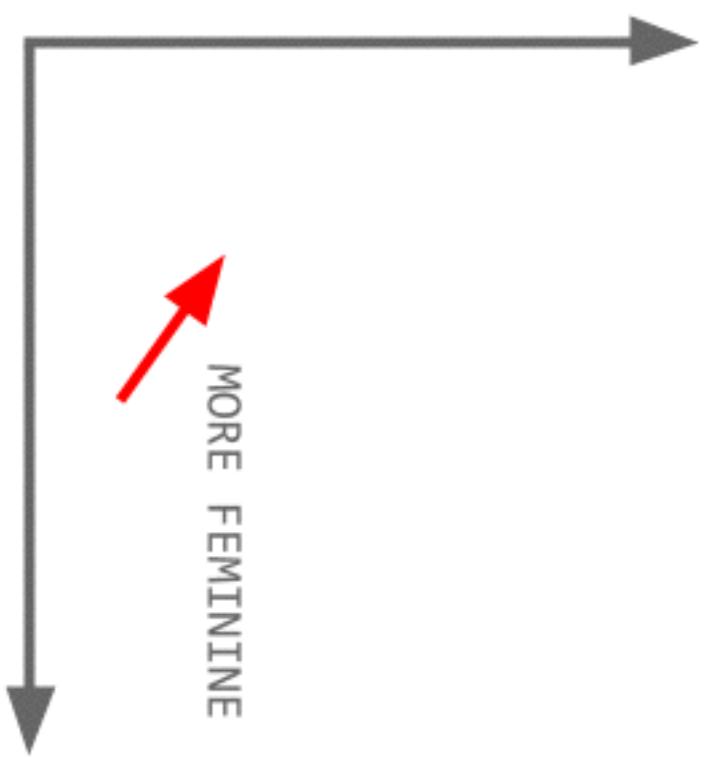
The red direction encodes gender



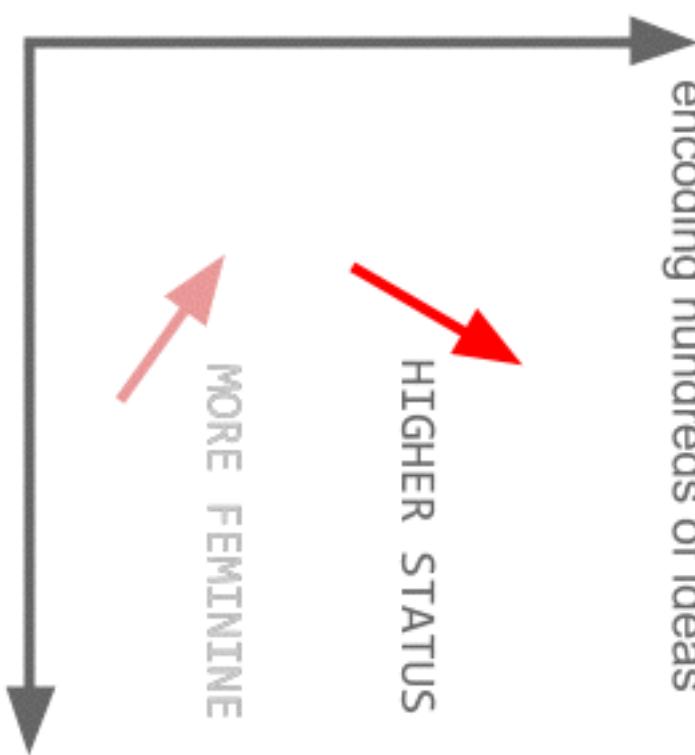
Which is consistent across all words



This **direction** always means **gender**



We have hundreds of **directions**  
encoding hundreds of ideas



Czech + currency	Vietnam + capital	German + airlines	Russian + river	French + actress
koruna	Hanoi	airline Lufthansa	Moscow	Juliette Binoche
Check crown	Ho Chi Minh City	carrier Lufthansa	Volga River	Vanessa Paradis
Polish zolty	Viet Nam	flag carrier Lufthansa	upriver	Charlotte Gainsbourg
CTK	Vietnamese	Lufthansa	Russia	Cecile De

ITEM\_3469 + 'Pregnant',



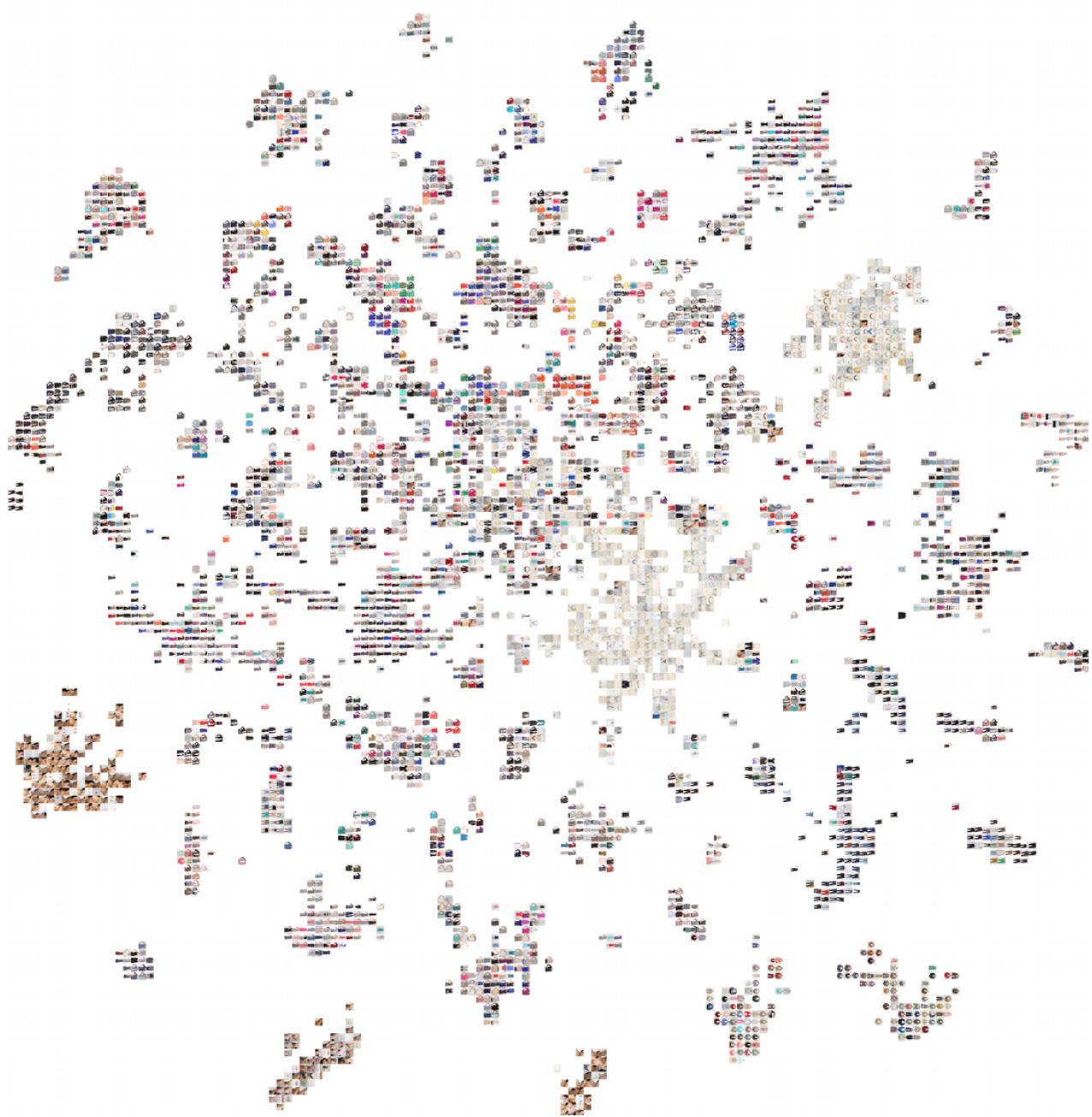
+ 'Pregnant'

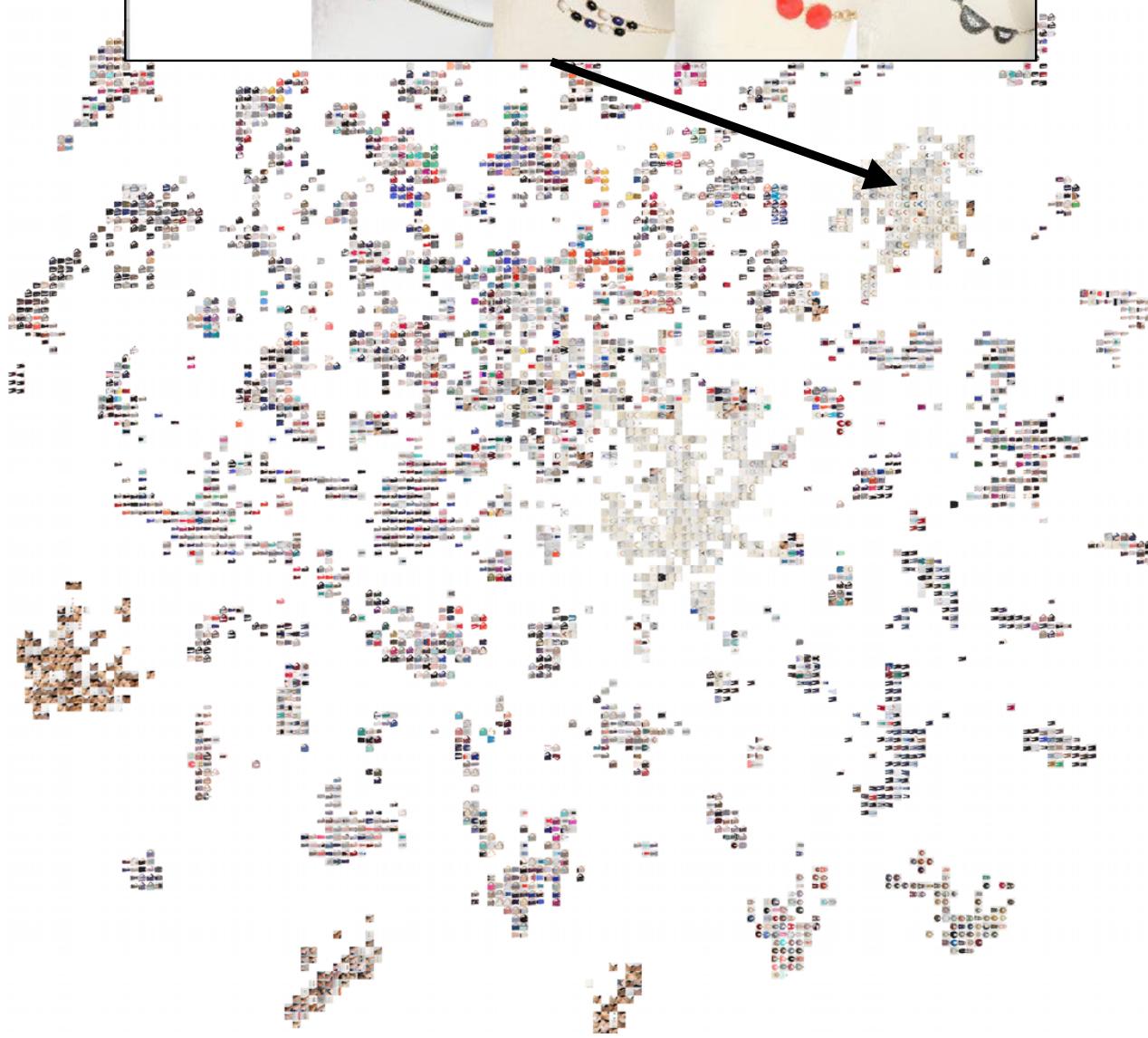
= ITEM\_701333  
= ITEM\_901004  
= ITEM\_800456

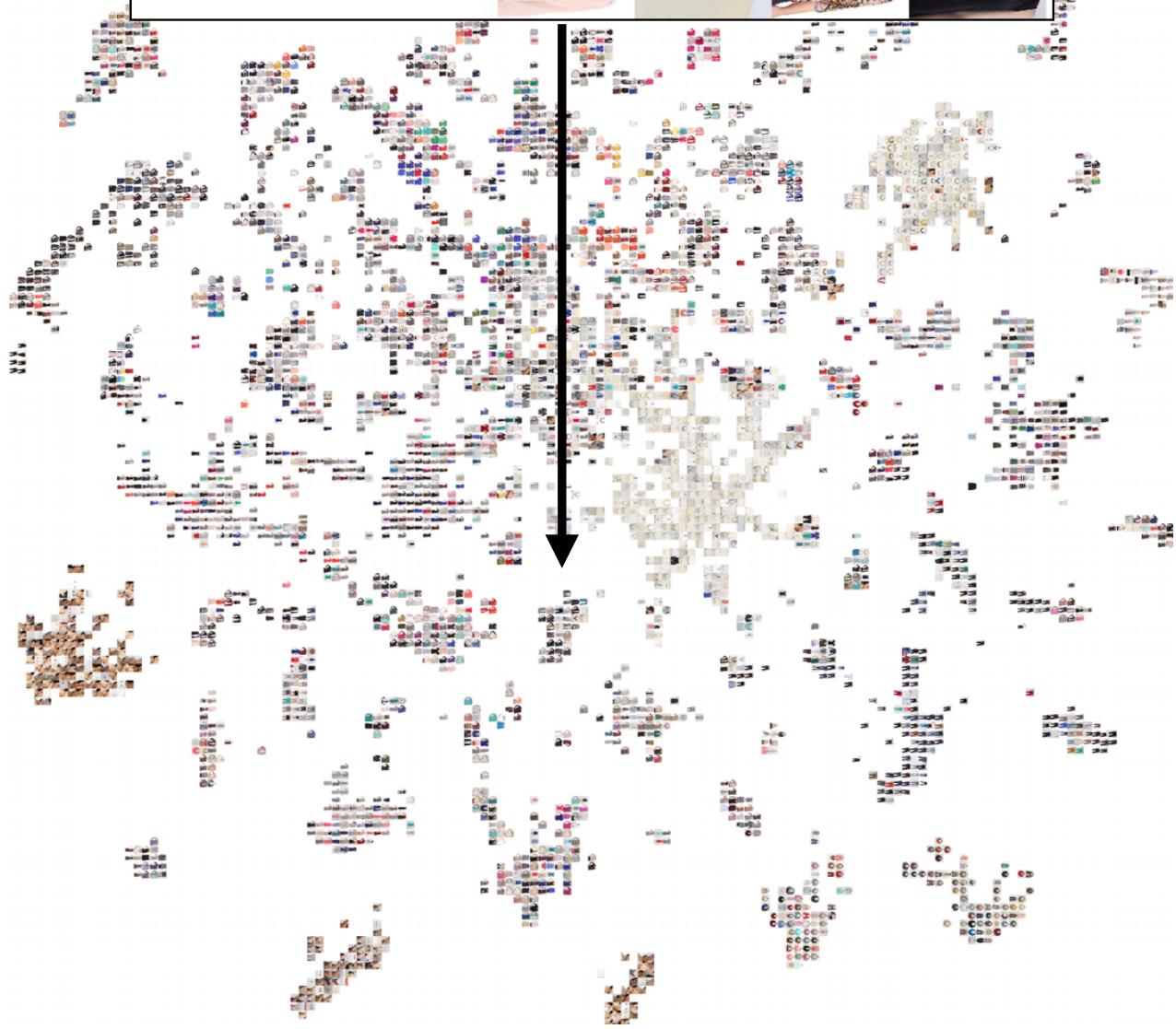
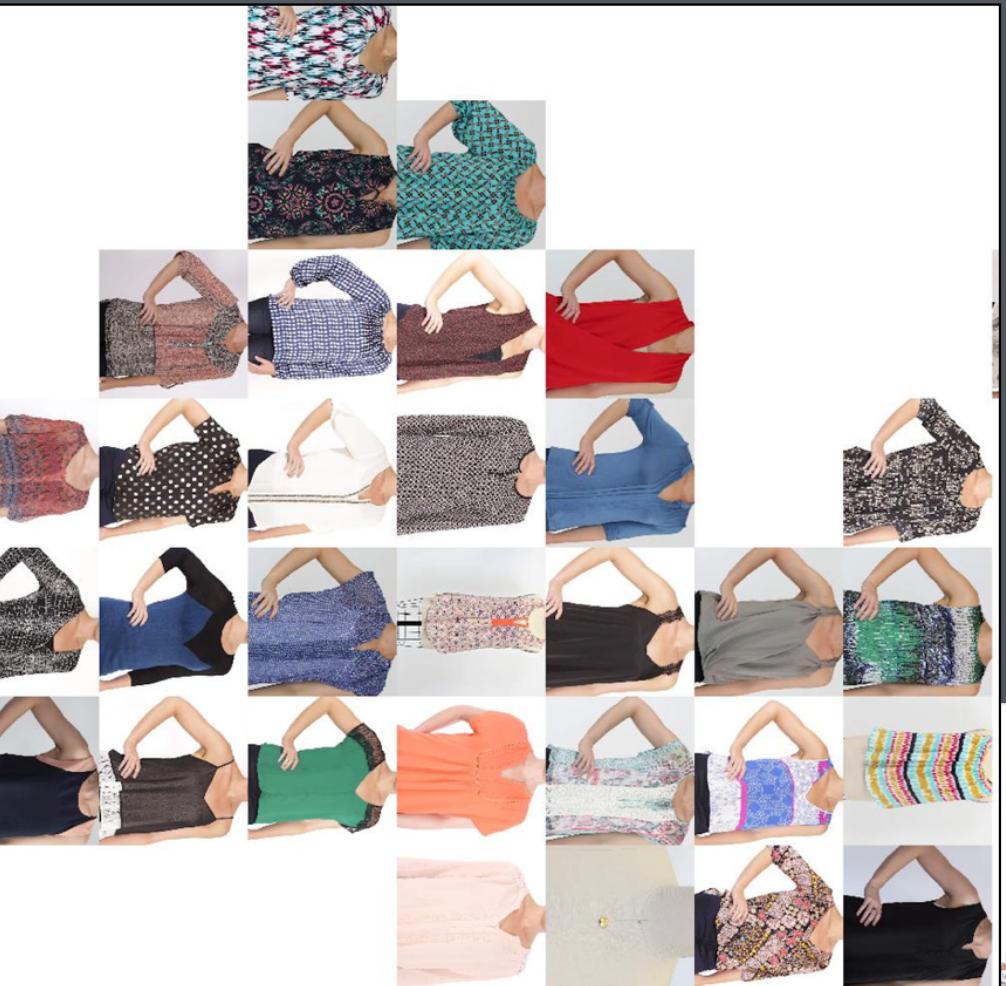


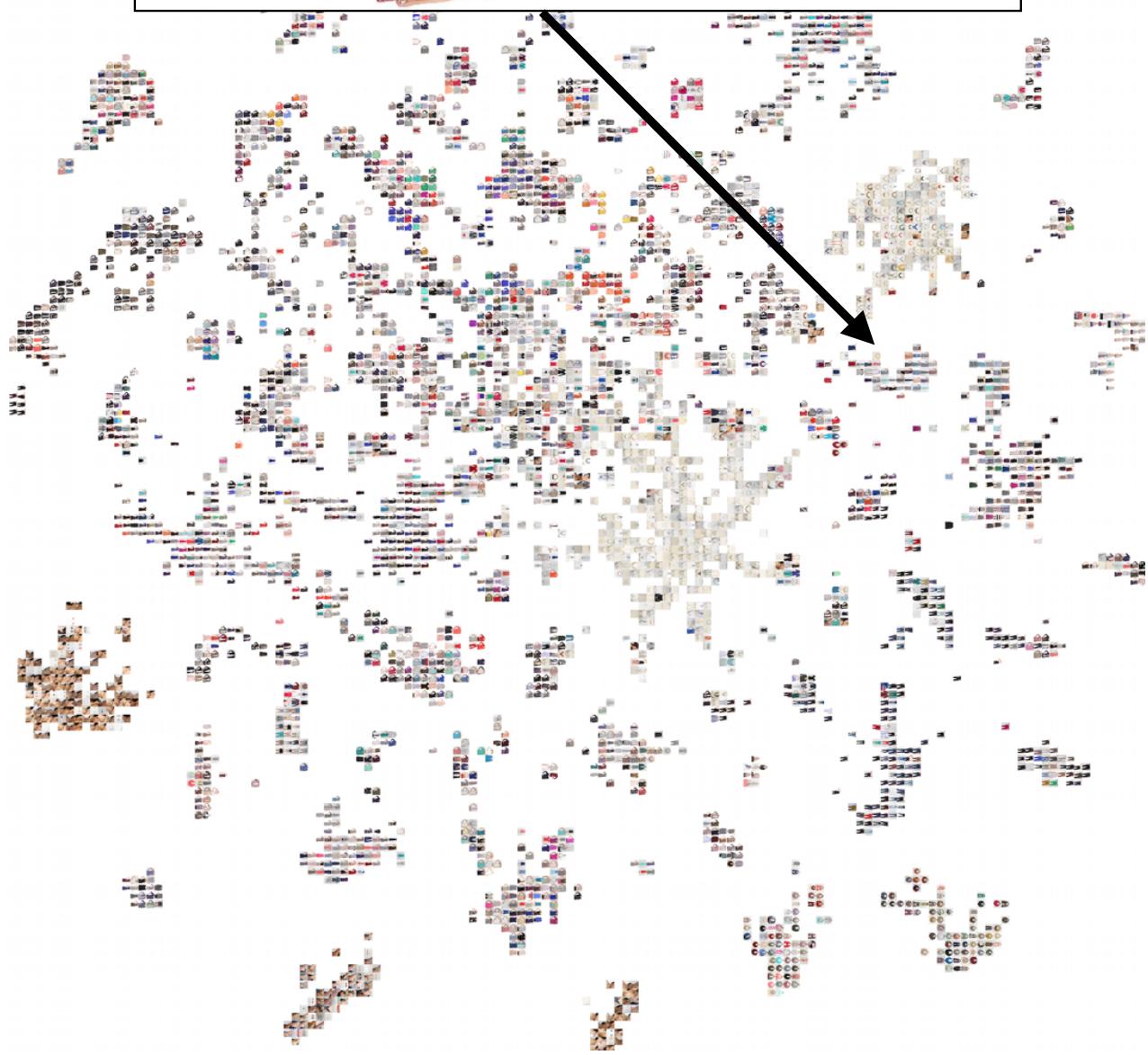
What about LDA?

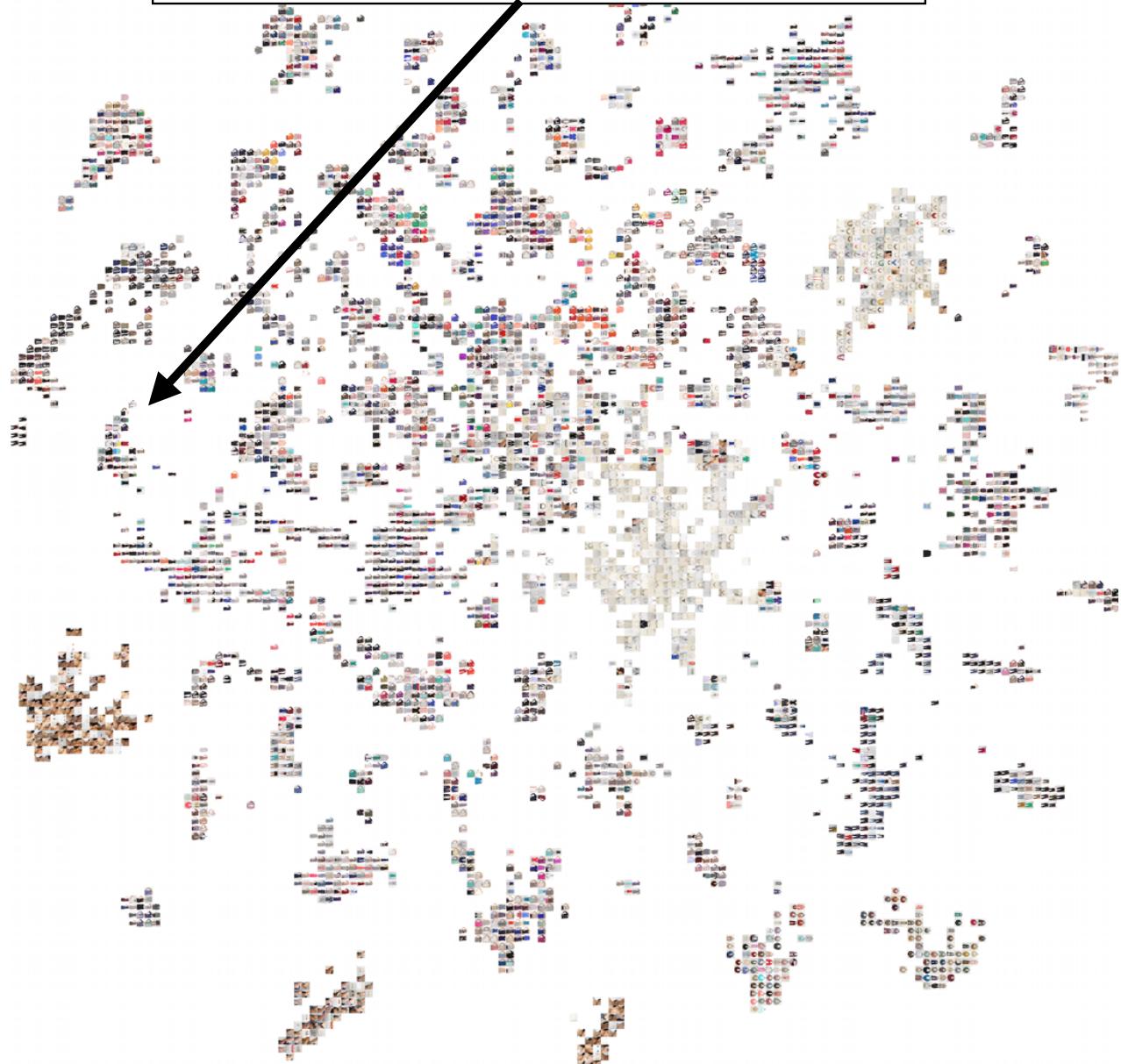
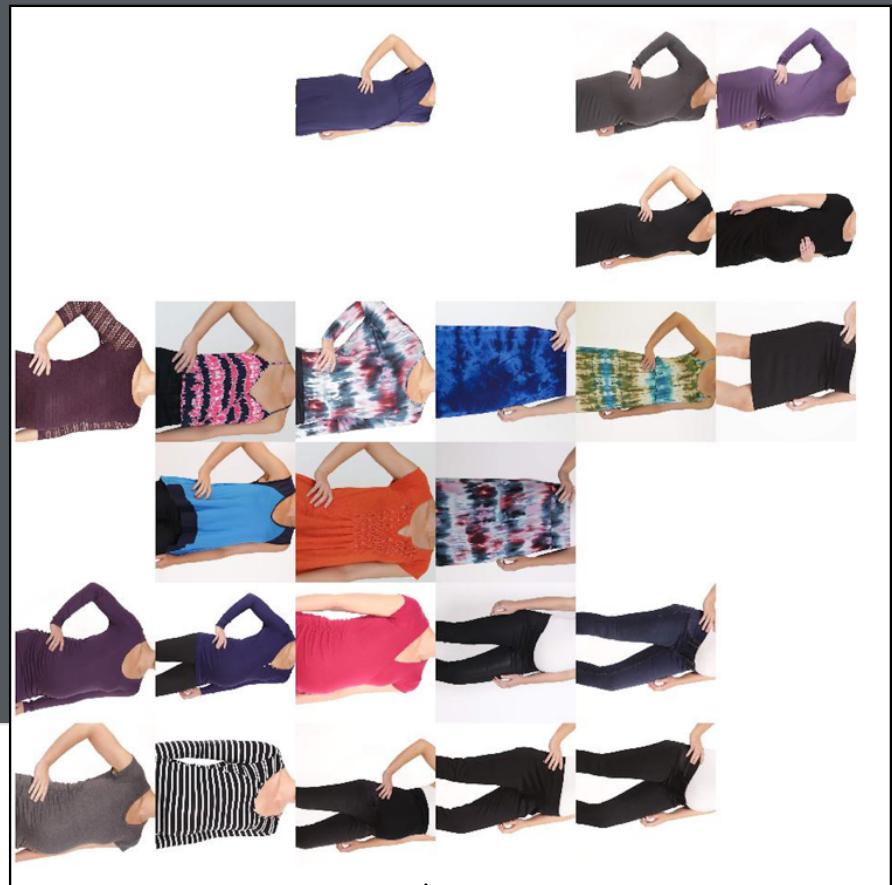
# LDA on Client Item Descriptions

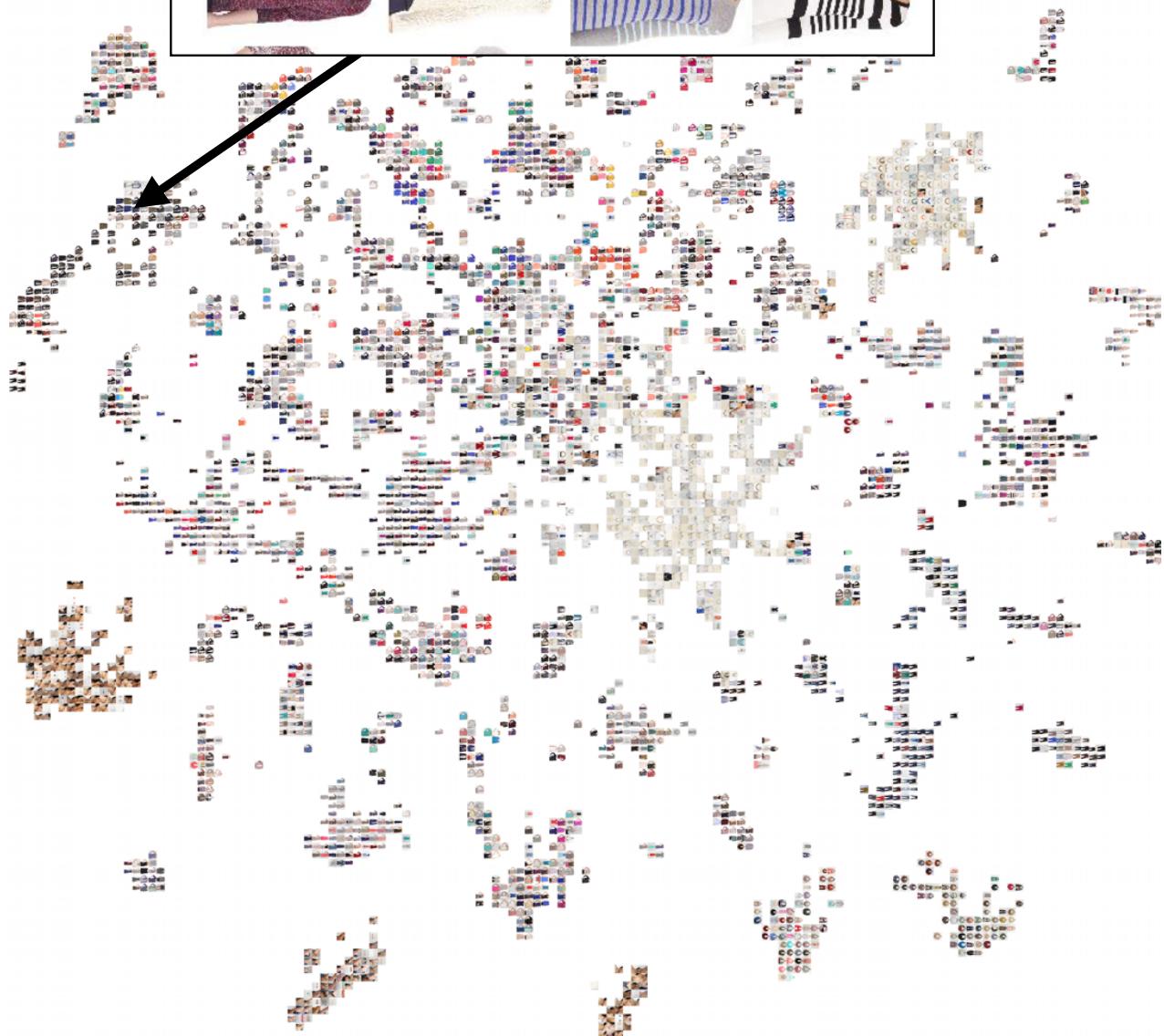




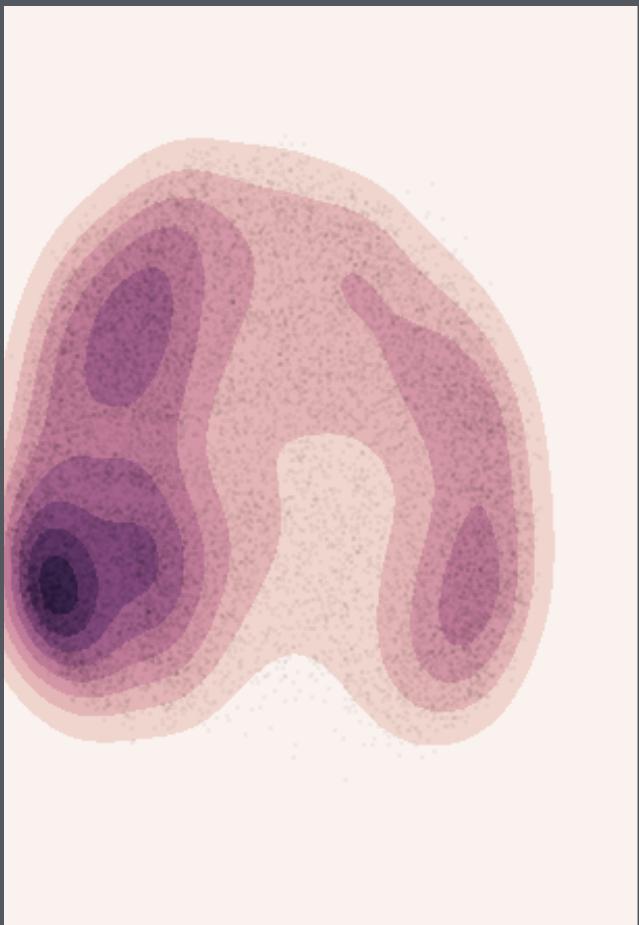




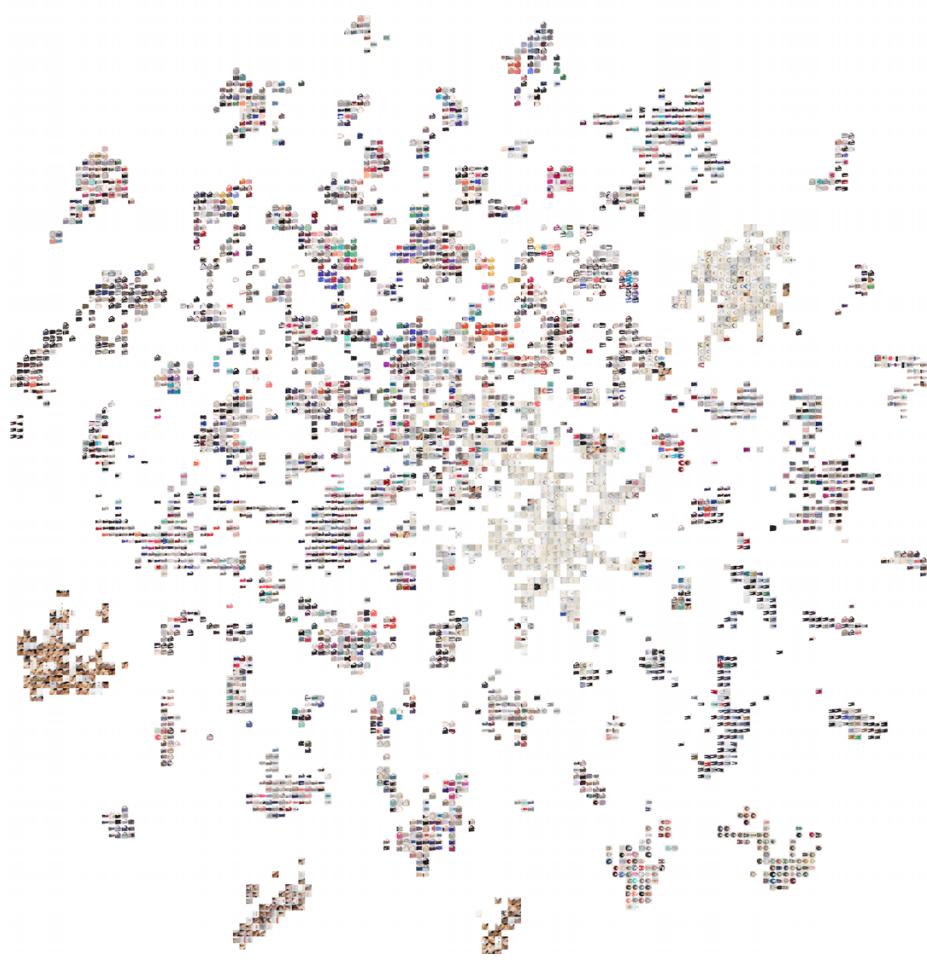




Pairwise gamma correlation  
from style ratings



Latent style vectors from text



Diversity from ratings

Diversity from text

lida vs word2vec

“I love finding new designer brands for jeans”



word2vec is *local*:  
one **word** predicts a nearby **word**

“I love finding new designer brands for jeans”

client\_comments

I really like the color of this top and the fit but for suc...

Almost too big. Love the dress though. Going to k...

EVERYTHING about this dress is absolutely PERFE...

This was a Winner to Update my look.... thanks...

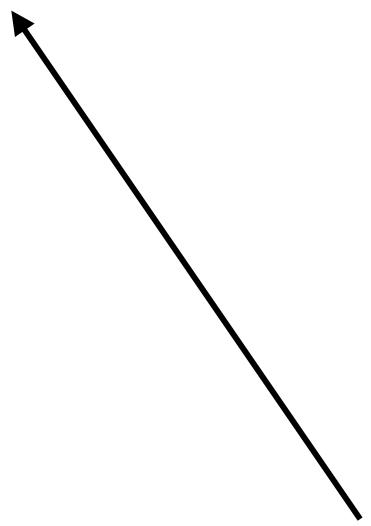
Love love love!!! Nothing more to say here.

I love finding new designer brands for jeans. I usuall...

Didn't think I'd be too interested in jewelry but t...

Love love love the color, pattern and flowiness!

But text is usually organized.



“I love finding new designer brands for jeans”

But text is usually organized.

client_comments	document_id
I really like the color of this top and the fit but for suc...	5943
Almost too big. Love the dress though. Going to k...	5872
EVERYTHING about this dress is absolutely PERFE...	5951
This was a Winner to Update my look.... thanks...	4017
Love love love!!! Nothing more to say here.	5953
I love finding new designer brands for jeans. I usuall...	7681
Didn't think I'd be too interested in jewelry but t...	3870
Love love love the color, pattern and flowiness!	6286

“I love finding new designer brands for jeans”

doc 7681



In LDA, documents *globally* predict *words*.

client_comments	document_id
I really like the color of this top and the fit but for suc...	5943
Almost too big. Love the dress though. Going to k...	5872
EVERYTHING about this dress is absolutely PERFE...	5951
This was a Winner to Update my look.... thanks...	4017
Love love love!!! Nothing more to say here.	5953
I love finding new designer brands for jeans. I usuall...	7681
Didn't think I'd be too interested in jewelry but t...	3870
Love love love the color, pattern and flowiness!	6286

typical word2vec vector

[ -0.75, -1.25, -0.55, -0.12, +2.2]

typical LDA document vector

[ 0%, 9%, 78%, 11%]

typical word2vec vector

[ -0.75, -1.25, -0.55, -0.12, +2.2]

typical LDA document vector

[ 0%, 9%, **78%**, 11%]

All real values

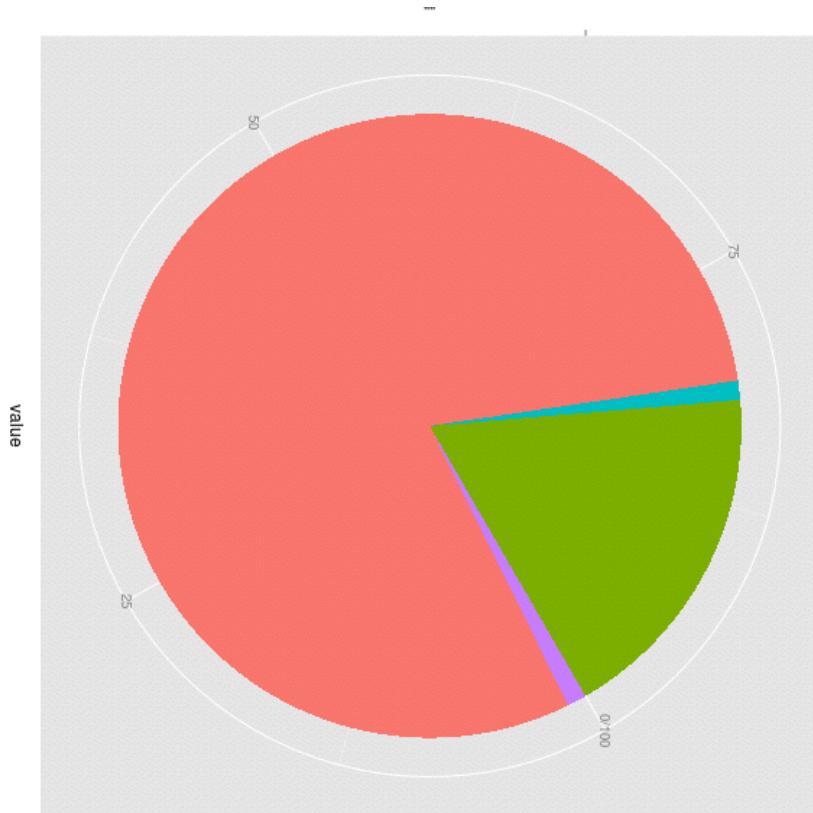
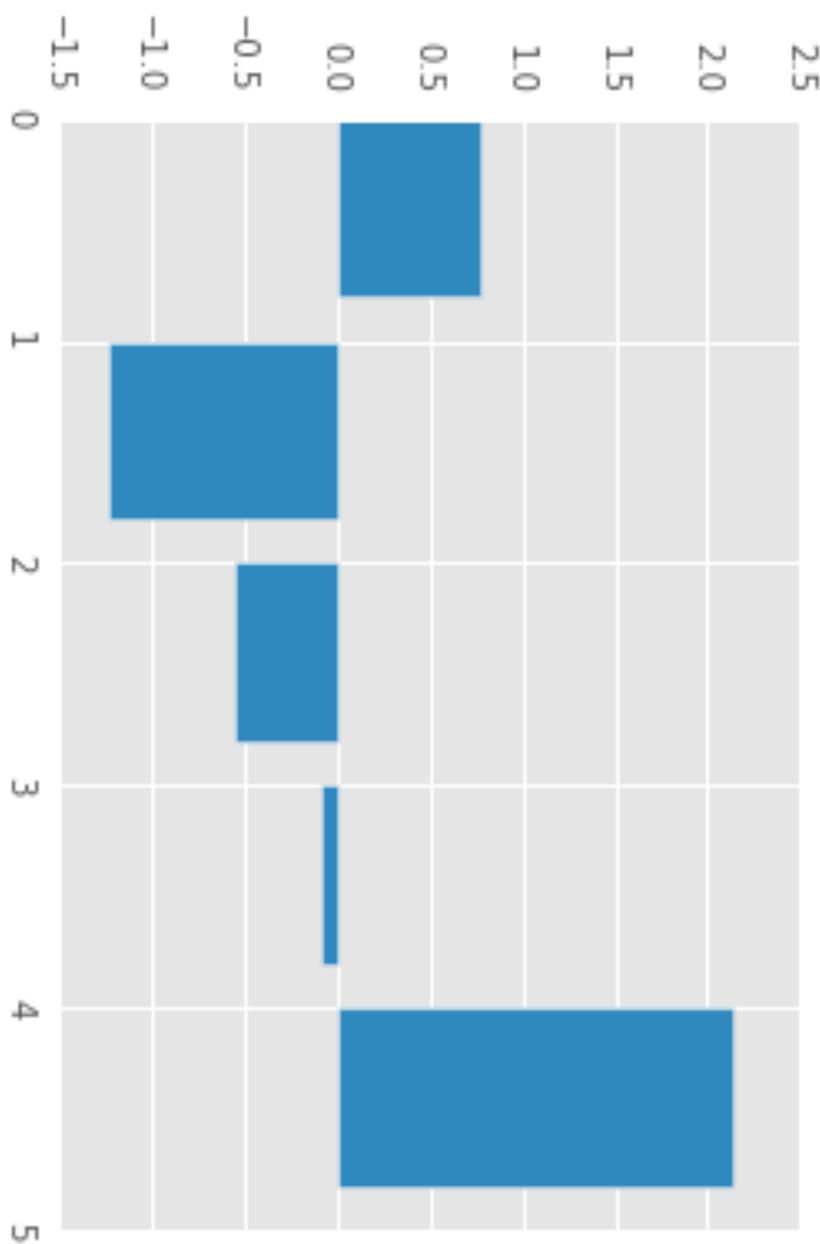
All sum to 100%

5D word2vec vector

[ -0.75, -1.25, -0.55, -0.12, +2.2]

5D LDA document vector

[ 0%, 9%, **78%**, 11%]



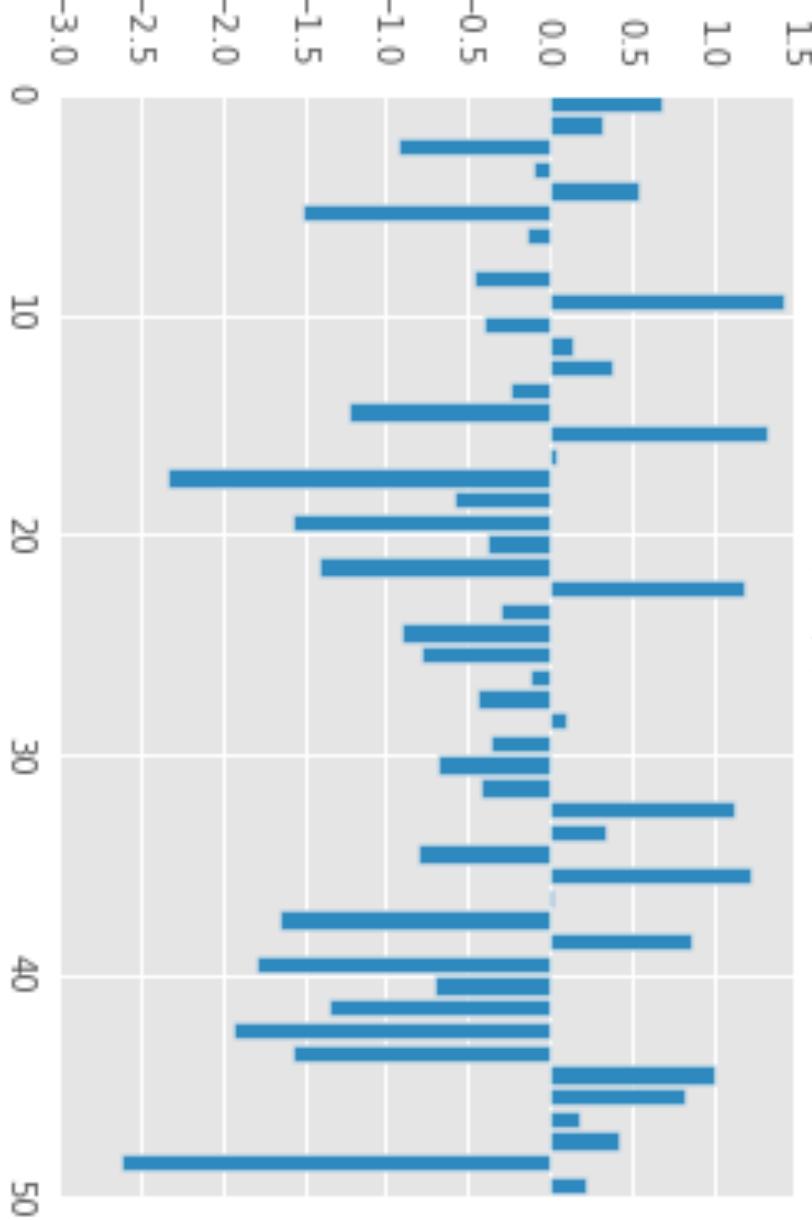
Topics  
Tops  
Jewelry  
Denim  
Bottoms

# 100D word2vec vector

[ -0.75, -1.25, -0.55, -0.27, -0.94, 0.44, 0.05, 0.31 ... -0.12, +2.2 ]

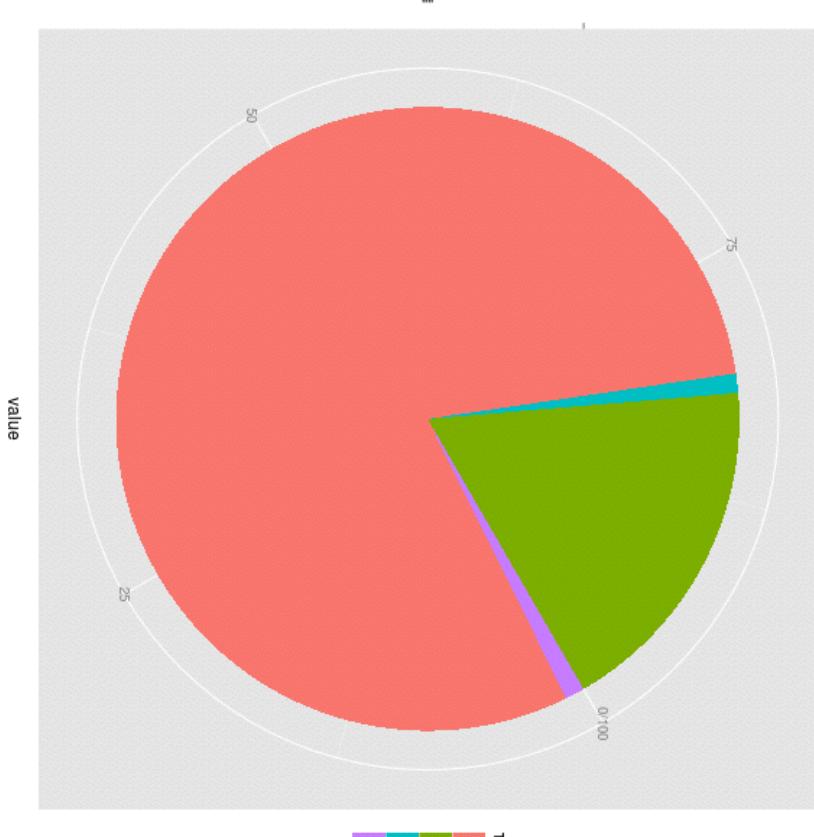
[ 0% 0% 0% 0% 0% ... 0%, 9%, 78%, 11% ]

**dense**



# 100D LDA document vector

**sparse**



100D word2vec vector

[ -0.75, -1.25, -0.55, -0.27, -0.94, 0.44, 0.05, 0.31 ... -0.12, +2.2 ]

[ 0% 0% 0% 0% 0% ... 0%, 9%, 78%, 11% ]

100D LDA document vector

Similar in 100D ways  
(very **flexible**)

+mixture  
+sparse

can we do both? **Idavec**

 @chrisemoody

The goal:  
Use all of this context to learn  
interpretable topics.

client\_comments

I love finding new designer  
brands for jeans. I usual...  
Didn't think I'd be too  
interested in jewelry but t...

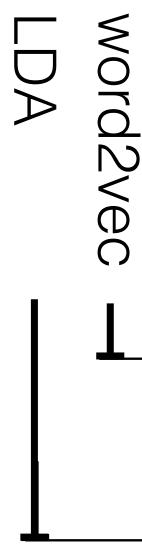
word2vec  $\rightarrow P(v_{OUT} | v_{IN})$

The goal:  
Use all of this context to learn  
interpretable topics.

client_comments	document_id
[REDACTED]	5943
[REDACTED]	5872
[REDACTED]	5951
[REDACTED]	4017
[REDACTED]	5953
I love finding new designer brands for jeans. I usually... Didn't think I'd be too interested in jewelry but t...	7681
[REDACTED]	3870
[REDACTED]	6286

this document is  
80% high fashion

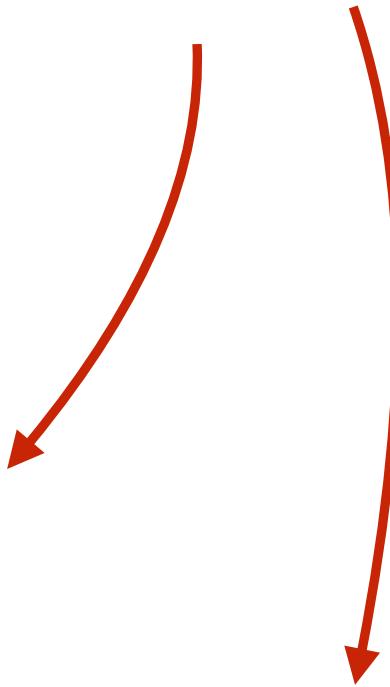
this document is  
60% style



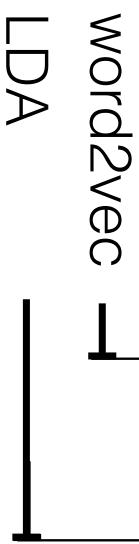
The goal:  
Use all of this context to learn  
interpretable topics.

client_comments	document_id	zip_code
[REDACTED]	5943	52
[REDACTED]	5872	194
[REDACTED]	5951	158
[REDACTED]	4017	991
[REDACTED]	5953	193
I love finding new designer brands for jeans. I usual...	7681	314
Didn't think I'd be too interested in jewelry but t...	3870	43
[REDACTED]	6286	151

this zip code is  
80% hot climate



this zip code is  
60% outdoors wear



The goal:  
Use all of this context to learn  
interpretable topics.

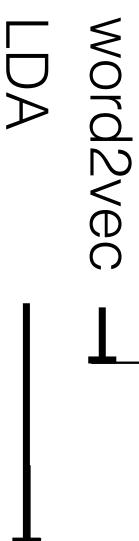
client_comments	document_id	zip_code	client_id
[REDACTED]	5943	52	5977
[REDACTED]	5872	194	5906
[REDACTED]	5951	158	5985
[REDACTED]	4017	991	4051
[REDACTED]	5953	193	5987
I love finding new designer brands for jeans. I usual...	7681	314	7715
Didn't think I'd be too interested in jewelry but t...	3870	43	3904
[REDACTED]	6286	151	6320

this client is  
80% sporty

60% casual wear

I love finding new designer brands for jeans. I usual...

Didn't think I'd be too interested in jewelry but t...



“PS! Thank you for such an awesome top”



word2vec predicts *locally*:

one **word** predicts a nearby **word**

$$P(v_{OUT} | v_{IN})$$

l<sub>a</sub>2vec

$v_{DOC}$

doc\_id=1846

“PS! Thank you for such an awesome top

$v_{OUT}$



LDA predicts a *word* from a *global* context

$P(v_{OUT} | v_{DOC})$

$v_{DOC}$ 

doc\_id=1846

“PS! Thank you for such an awesome top

 $v_{IN}$      $v_{OUT}$ 

can we predict a **word** both *locally* and *globally*?

$v_{DOC}$ 

doc\_id=1846

“PS! Thank you for such an awesome top

 $v_{IN}$      $v_{OUT}$ 

can we predict a **word** both *locally* and *globally*?

$$P(v_{OUT} | v_{IN} + v_{DOC})$$

$v_{DOC}$

$v_{IN} \quad v_{OUT}$ ,  
doc\_id=1846 “PS! Thank you for such an awesome top”



can we predict a **word** both *locally* and *globally*?

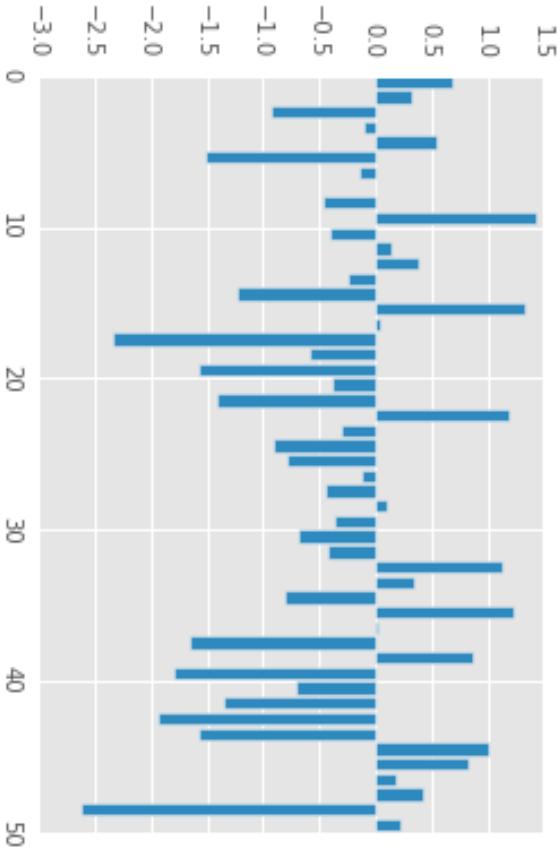
$$P(v_{OUT} | v_{IN} + v_{DOC})$$



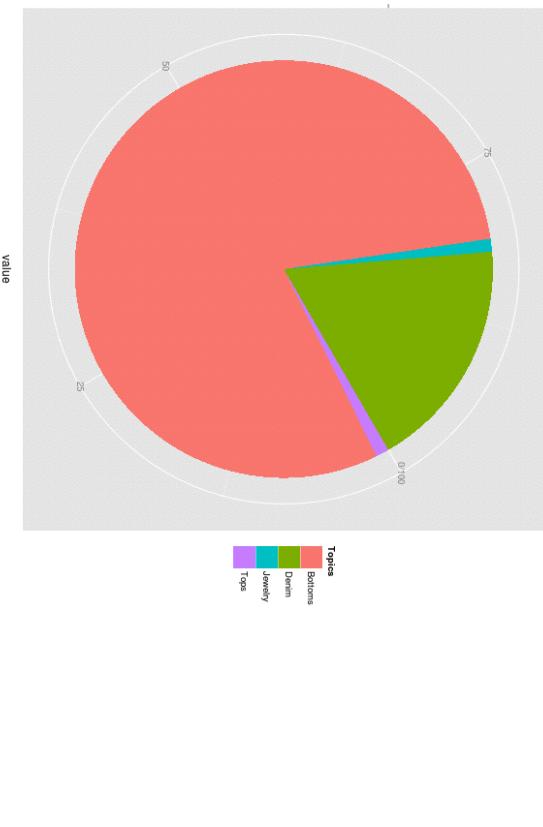
\*very similar to the Paragraph Vectors / doc2vec

This works! 😊 But *v<sub>doc</sub>* isn't as interpretable as the LDA topic vectors. 😞

This works! 😊 But *v<sub>DOC</sub>* isn't as interpretable as the LDA topic vectors. 😞



This works! 😊 But *v\_doc* isn't as interpretable as the LDA topic vectors. 😞



This works! 😊 But *v<sub>doc</sub>* isn't as

interpretable as the LDA topic vectors. 😞

We're missing *mixtures & sparsity*.

This works! 😊 But  $v_{DOC}$  isn't as

interpretable as the LDA topic vectors. 🙁

Let's make  $v_{DOC}$  into a mixture...

Let's make  $v_{DOC}$  into a mixture...

$$v_{DOC} = \alpha v_{topic1} + \beta v_{topic2} + \dots$$

(up to k topics)

Let's make  $v_{DOC}$  into a mixture...

$$v_{DOC} = \alpha v_{topic1} + \beta v_{topic2} + \dots$$



*Trinitarian  
baptismal  
Pentecostals  
Bede  
schismatics  
excommunication*

Let's make  $v_{DOC}$  into a mixture...

topic 1 = “religion”,

*Trinitarian  
baptismal*

$$v_{DOC} = \alpha v_{topic1} + \beta v_{topic2} + \dots$$



*Pentecostals*

*Bede*

*schismatics*

*excommunication*

Let's make  $v_{DOC}$  into a mixture...

**topic 1 = “religion”**

*Trinitarian  
baptismal*

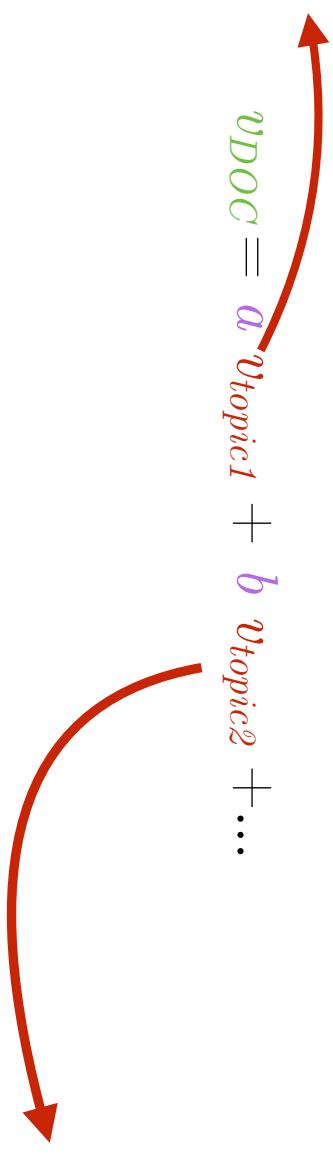
$$v_{DOC} = \alpha v_{topic1} + \beta v_{topic2} + \dots$$

*Pentecostals*

*Bede*

*schismatics*

*excommunication*



*Milosevic*

*absentee*

*Indonesia*

*Lebanese*

*Isrealis*

*Karadzic*

Let's make  $v_{DOC}$  into a mixture...

**topic 1 = “religion”**

*Trinitarian*

*baptismal*

*Pentecostals*

*bede*

*schismatics*

*excommunication*

$$v_{DOC} = \alpha v_{topic1} + \beta v_{topic2} + \dots$$

**topic 2 = “politics”**

*Milosevic*

*absentee*

*Indonesia*

*Lebanese*

*Isrealis*

*Karadzic*

Let's make  $v_{DOC}$  into a mixture...

topic 1 = “religion”

*Milosevic*

*absentee*

*Indonesia*

*Lebanese*

*Isrealis*

*Karadzic*

$$v_{DOC} = 10\% \text{ religion} + 89\% \text{ politics} + \dots$$



topic 2 = “politics”

*Trinitarian*

*baptismal*

*Pentecostals*

*bede*

*schismatics*

*excommunication*



Let's make  $v_{DOC}$  sparse

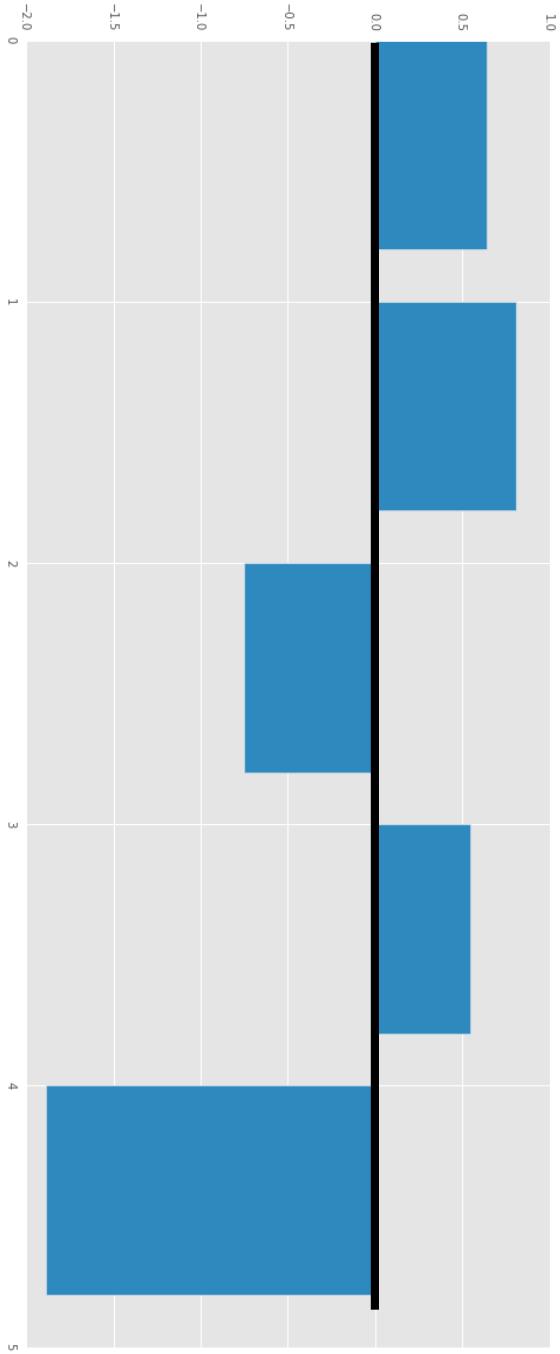
$$v_{DOC} = \alpha v_{religion} + \beta v_{politics} + \dots$$



$$[ -0.75, \quad -1.25, \quad \dots ]$$

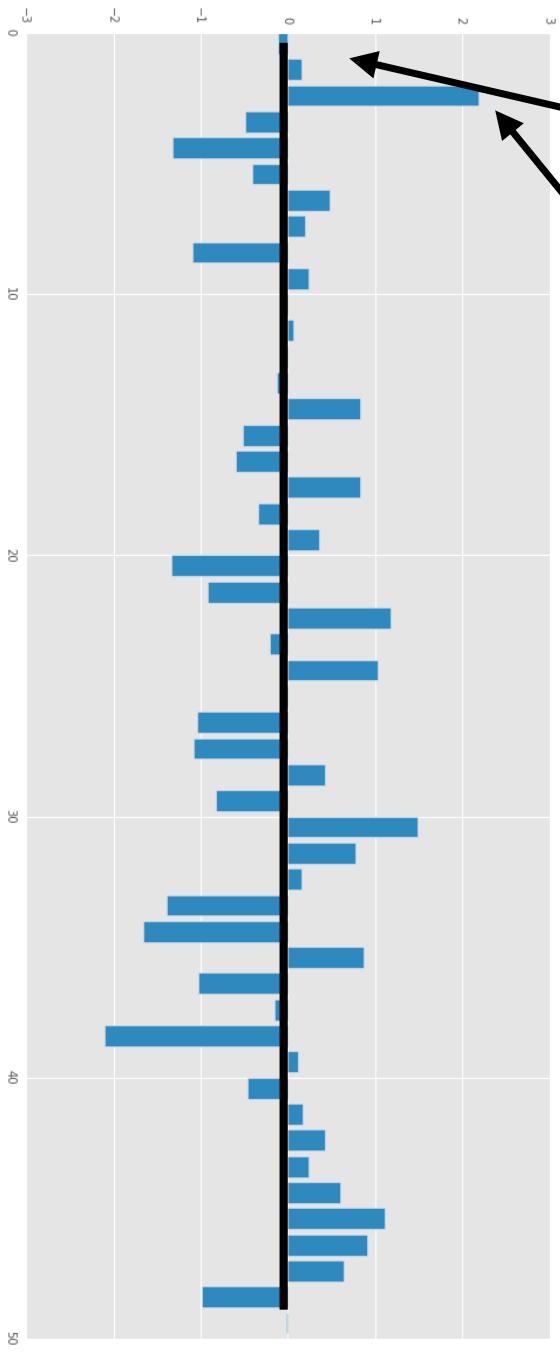
Let's make  $v_{DOC}$  sparse

$$v_{DOC} = \alpha v_{religion} + \beta v_{politics} + \dots$$



Let's make  $v_{DOC}$  sparse

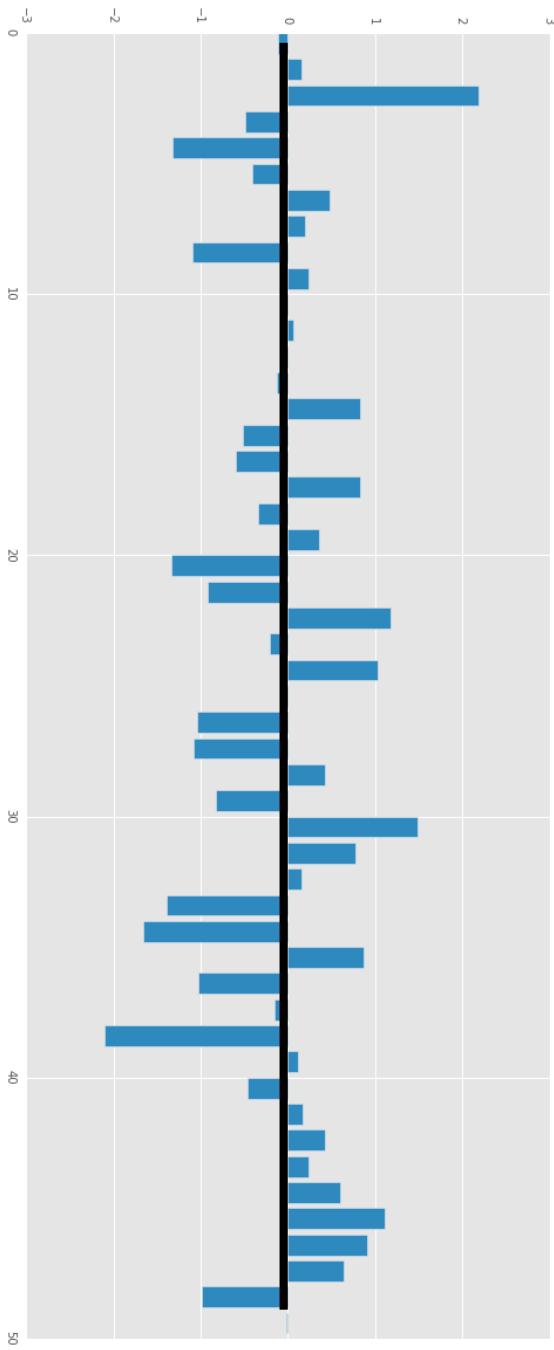
$$v_{DOC} = \alpha v_{religion} + \beta v_{politics} + \dots$$



Let's make  $v_{DOC}$  sparse

$$v_{DOC} = \alpha v_{religion} + \beta v_{politics} + \dots$$

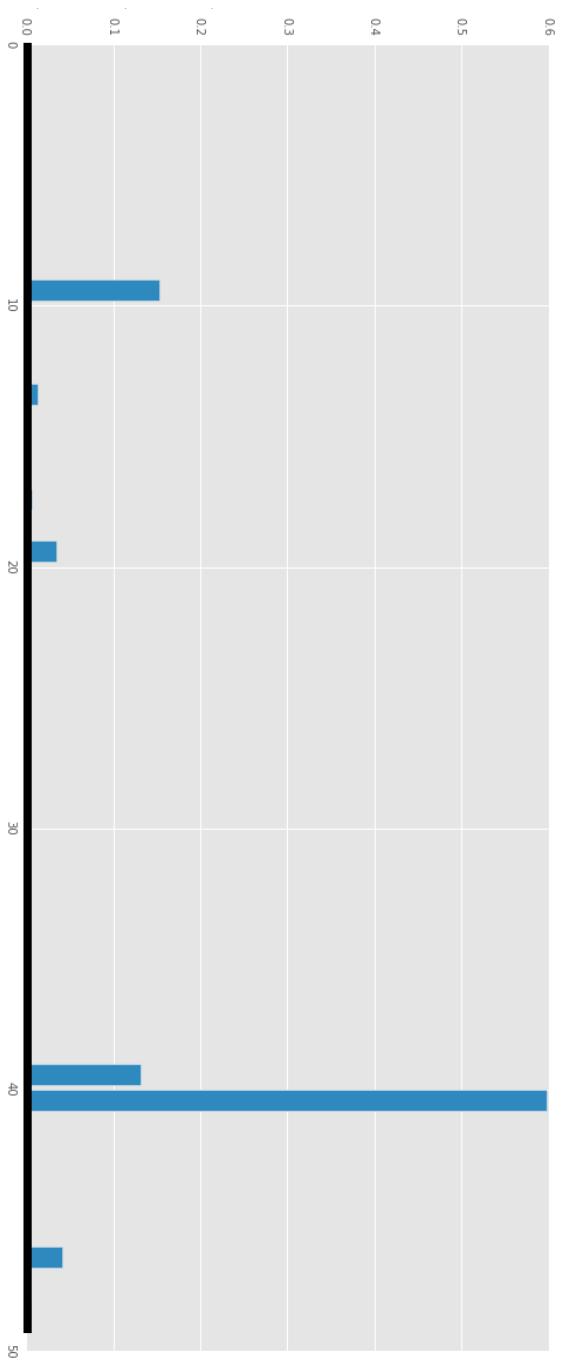
$$\{a, b, c, \dots\} \sim dirichlet(\alpha)$$



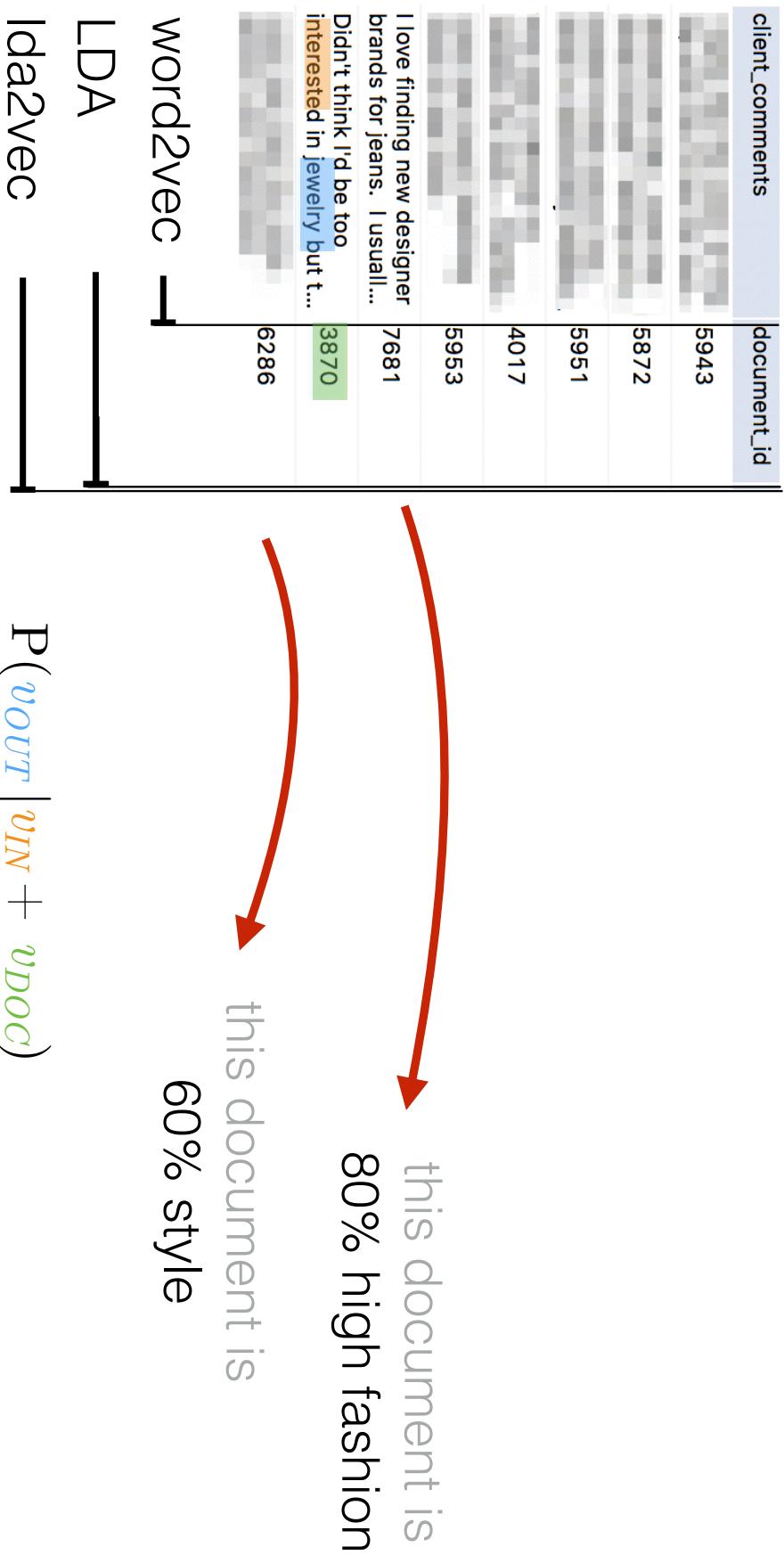
Let's make  $v_{DOC}$  sparse

$$v_{DOC} = \alpha v_{religion} + \beta v_{politics} + \dots$$

$$\{a, b, c, \dots\} \sim dirichlet(\alpha)$$



The goal:  
Use all of this context to learn  
interpretable topics.



The goal:  
 Use all of this context to learn  
 interpretable topics.

client_comments	document_id	zip_code
[REDACTED]	5943	52
[REDACTED]	5872	194
[REDACTED]	5951	158
[REDACTED]	4017	991
I love finding new designer brands for jeans. I usual...	5953	193
Didn't think I'd be too interested in jewelry but t...	7681	314
[REDACTED]	3870	43
[REDACTED]	6286	151

$$P(v_{OUT} | v_{IN} + v_{DOC} + v_{ZIP})$$

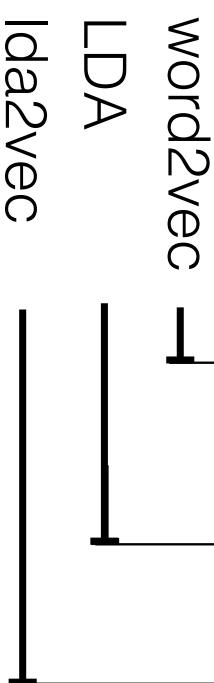
Word2Vec  
 LDA  
 lda2vec

The goal:  
Use all of this context to learn  
interpretable topics.

client_comments	document_id	zip_code
[REDACTED]	5943	52
[REDACTED]	5872	194
[REDACTED]	5951	158
[REDACTED]	4017	991
[REDACTED]	5953	193
I love finding new designer brands for jeans. I usual...	7681	314
Didn't think I'd be too interested in jewelry but t...	3870	43
[REDACTED]	6286	151

this zip code is  
80% hot climate

this zip code is  
60% outdoors wear



$$P(v_{OUT} | v_{IN} + v_{DOC} + v_{ZIP})$$

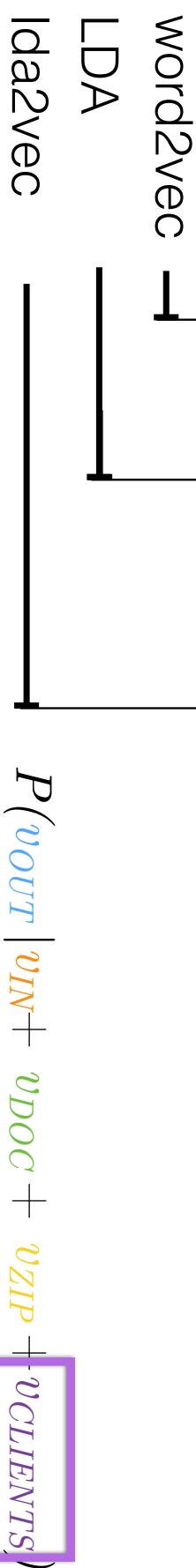
The goal:  
Use all of this context to learn  
interpretable topics.

client_comments	document_id	zip_code	client_id
[REDACTED]	5943	52	5977
[REDACTED]	5872	194	5906
[REDACTED]	5951	158	5985
[REDACTED]	4017	991	4051
[REDACTED]	5953	193	5987
I love finding new designer brands for jeans. I usual...	7681	314	7715
Didn't think I'd be too interested in jewelry but t...	3870	43	3904
[REDACTED]	6286	151	6320

I love finding new designer  
brands for jeans. I usual...  
Didn't think I'd be too  
interested in jewelry but t...

this client is  
80% sporty

60% casual wear



 @chrisemoody

The goal:  
Use all of this context to learn  
interpretable topics.

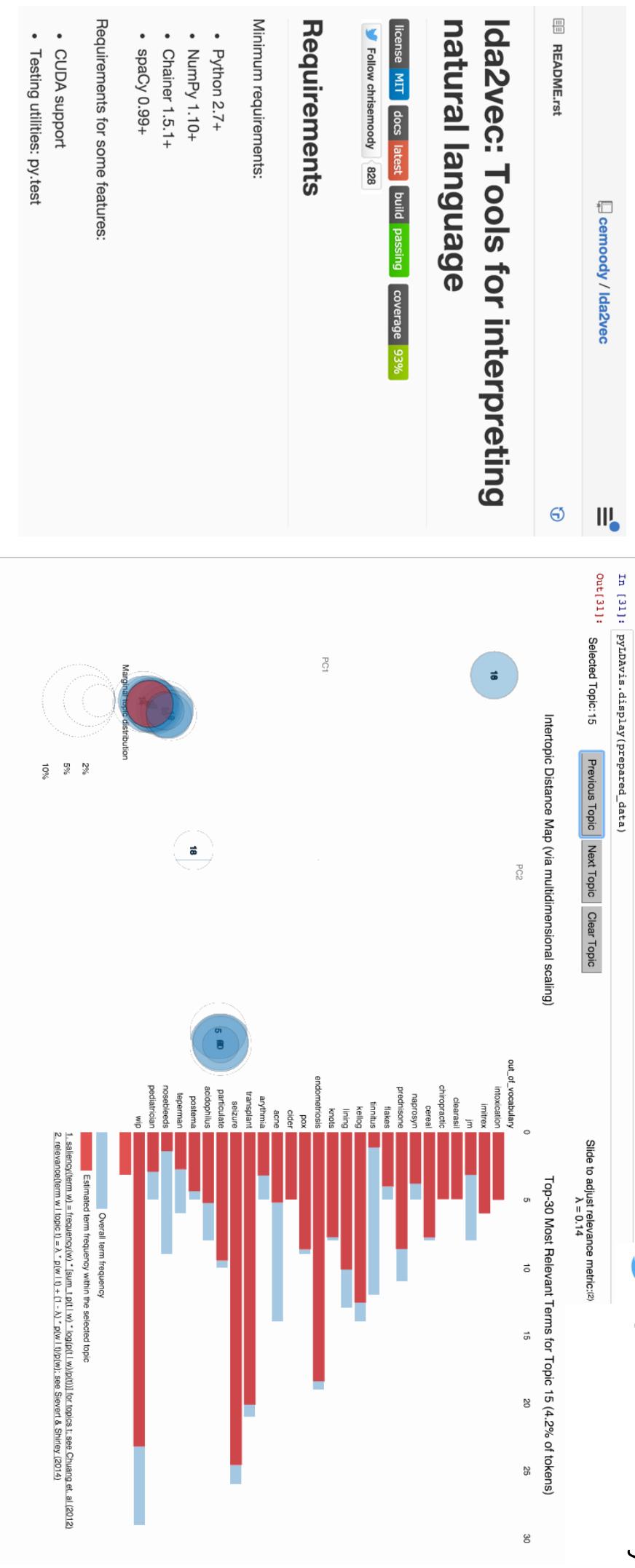
Can also make the topics *supervised* so that they predict an outcome.

The figure illustrates the relationship between client comments and document IDs. A red arrow points from the highlighted comment "I love finding new designer brands for jeans. I usually... Didn't think I'd be too interested in jewelry but..." to the document ID 3870. Another red arrow points from document ID 3870 to the document ID 3904.

client_comments	document_id	zip_code	client_id	sold
[REDACTED]	5943	52	5977	1
[REDACTED]	5872	194	5906	1
[REDACTED]	5951	158	5985	1
[REDACTED]	4017	991	4051	1
[REDACTED]	5953	193	5987	1
I love finding new designer brands for jeans. I usually... Didn't think I'd be too interested in jewelry but...	3870	314	7715	1
[REDACTED]	43	3904	1	1
[REDACTED]	6286	151	6320	1

$$P(\textcolor{blue}{vOUT} | \textcolor{orange}{vIN} + vDOC + \textcolor{brown}{vZIP} + vCLIENTS)$$

 @chrisemoody



API Ref docs (no narrative docs)  
GPU  
Decent test coverage

[github.com/chrisemoody/lda2vec](https://github.com/chrisemoody/lda2vec)

uses pyldavis

Can we model topics to sentences?

Ida2lstm

doc\_id=1846 “PS! Thank you for such an awesome idea”



@chrisemoody

doc\_id=1846  
“PS! Thank you for such an awesome idea”



Can we represent the internal LSTM  
states as a dirichlet mixture?

Can we model topics to sentences?

lida2lstm

doc\_id=1846  
“PS! Thank you for such an awesome idea”



Can we model topics to images?

lida2ae



TJ Torres



?



@chrisemoody

Multithreaded  
Stitch Fix





Bonus slides





## Paragraph Vectors

(Just extend the context window)

## Content dependency

(Change the window grammatically)

## Social word2vec (deepwalk)

(Sentence is a walk on the graph)

## Spotify

(Sentence is a playlist of song\_ids)

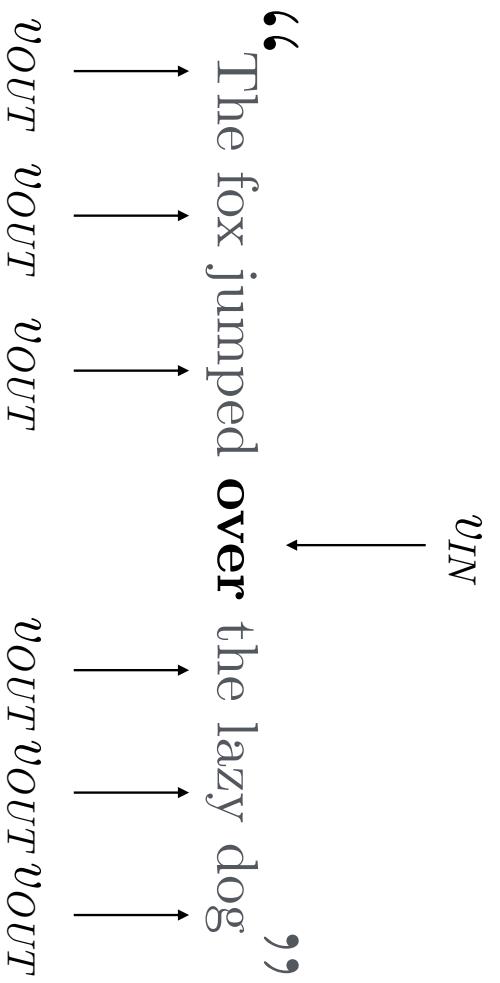
## Stitch Fix

(Sentence is a shipment of five items)

Relationship	Example 1	Example 2	Example 3
France - Paris big - bigger	Italy: Rome small: larger	Japan: Tokyo cold: colder	Florida: Tallahassee quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France copper - Cu	Berlusconi: Italy zinc: Zn	Merkel: Germany gold: Au	Koizumi: Japan uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

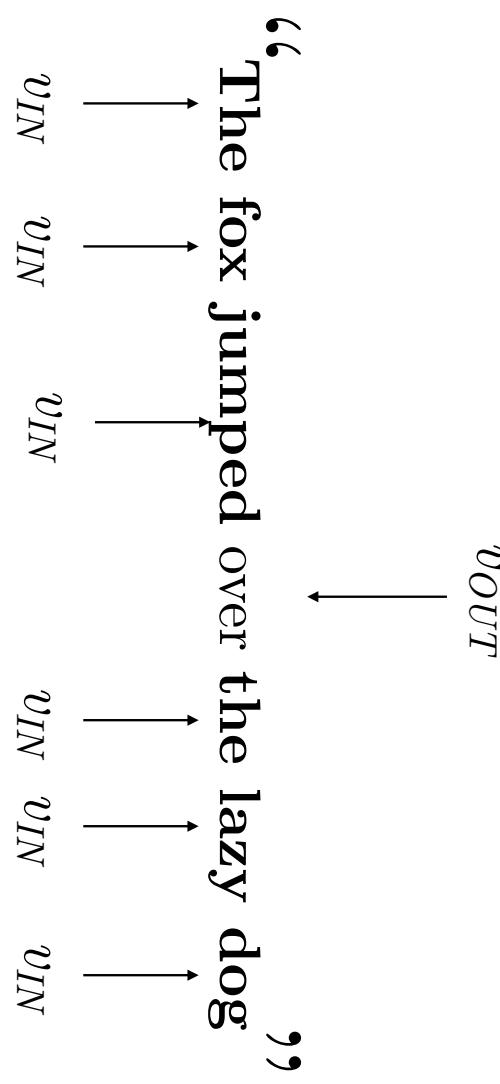
# SkipGram

Guess the context  
given the word



# CBOW

Guess the word  
given the context



Better at syntax.  
  
 $\sim 20x$  faster.

(this is the one we went over)

(this is the alternative.)



# LDA Results

Great Stylist

Perfect

I loved every choice in this fix!! Great job!

History

# LDA Results

Body Fit

My measurements are 36-28-32. If that helps.  
I like wearing some clothing that is fitted.  
Very hard for me to find pants that fit right.

# LDA Results

Sizing

Excited for next

Really enjoyed the experience and the pieces, sizing for tops was too big.  
Looking forward to my next box!

# LDA Results

Almost Bought

Perfect

It was a great fix. Loved the two items I kept and the three I sent back were close!

# What I didn't mention

---

A lot of text (only if you have a specialized vocabulary)

Cleaning the text

Memory & performance

Traditional databases aren't well-suited

False positives

and now for something **completely crazy**

All of the following ideas will change what  
'words' and 'context' represent.

What about summarizing documents?

On the day he took office, President Obama reached out to America's enemies, offering in his first inaugural address to **extend** a hand if you are willing to unclench your fist. More than six years later, he has arrived at a moment of truth in testing that

IN

On the day he took office, President Obama reached out to America's enemies, offering in his first inaugural address to **extend** a hand if you are willing to unclench your fist. More than six years later, he has arrived at a moment of truth in testing that

The framework nuclear agreement he reached with Iran on Thursday did not provide the definitive answer OUT whether Mr. Obama's audacious game OUT will pay off. The fist Iran has shaken at the so-called Great Satan since 1979 has not completely relaxed.

Normal skipgram extends  $C$  words before, and  $C$  words after.

IN

**doc\_1347**

OUT



On the day he took office, President Obama reached out to America's enemies, offering in his first inaugural address to extend a hand if you are willing to unclench your fist. More than six years later, he has arrived at a moment of truth in testing that

The framework nuclear agreement he reached with Iran on Thursday did not provide the definitive answer OUT whether Mr. Obama's audacious game OUT will pay off. The fist Iran has shaken at the so-called Great Satan since 1979 has not completely relaxed.

A document vector simply extends the context to the whole document.

```
from gensim.models import Doc2Vec
fn = "item_document_vectors"
model = Doc2Vec.load(fn)
model.most_similar('pregnant')
matches = list(filter(lambda x: 'SENT' in x[0], matches))

# [':...I am currently 23 weeks pregnant...',  
# ':...I'm now 10 weeks pregnant...',  
# ':...not showing too much yet...',  
# ':...15 weeks now. Baby bump...',  
# ':...6 weeks post partum!...',  
# ':...12 weeks postpartum and am nursing...',  
# ':...I have my baby shower that...',  
# ':...am still breastfeeding...',  
# ':...I would love an outfit for a baby shower...' ]
```

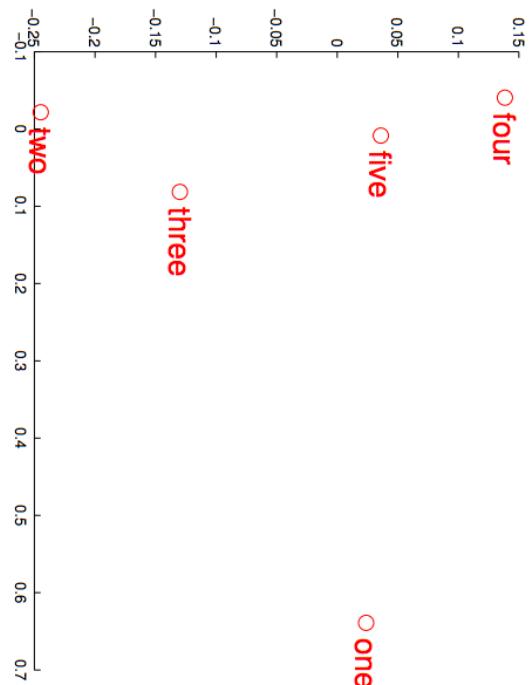
# translation

(using just a rotation  
matrix)

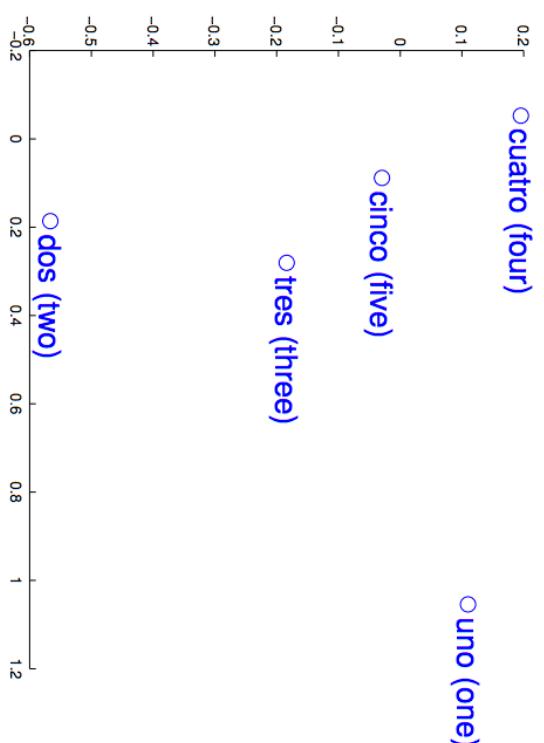
## Rotation Matrix



English



Spanish



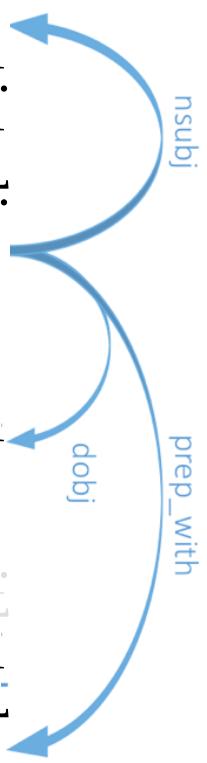
# context dependent

Australian scientist **discovers** star with telescope

context +/- 2 words

context  
dependent

Australian scientist discovers star with telescope



```
graph TD; Australian[Australian] -- "nsubj" --> NP1(( )); discovers[discover] -- "prep_with" --> NP2(( )); telescope[telescope] -- "dobj" --> NP3(( ));
```

# context dependent

Australian scientist discovers star with telescope

The diagram illustrates a context-dependent relationship. A horizontal double-headed arrow connects the words 'Australian scientist' and 'telescope'. Above this arrow, a curved blue arrow labeled 'prep\_with' points from 'Australian scientist' to 'telescope'. Below the arrow, the word 'context' is written vertically. To the left of the main text, there is a red diagonal banner containing the text 'Levy & Goldberg 2014'.

# context dependent

BoW	DEPS
dumbledore	sunnydale
hallows	collinwood
half-blood	calarts
malfoy	greendale
snape	millfield

topically-similar      vs      ‘functionally’ similar

Also show that SGNS is simply factorizing:

context  
dependent

$$w * c = PMI(w, c) - \log k$$

This is **completely** amazing!

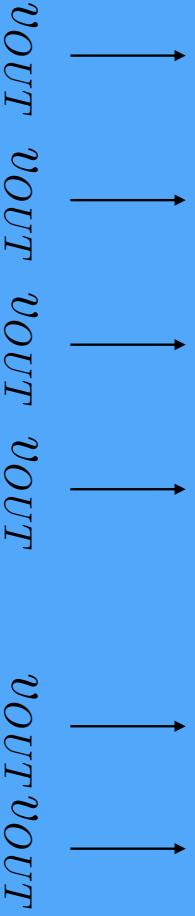
Intuition: positive associations (canada, snow)  
stronger in humans than negative associations  
(what is the opposite of Canada?)

## word2vec

learn word vectors from sentences

‘words’ are graph vertices  
‘sentences’ are random walks on the graph

‘The fox jumped over the lazy dog’

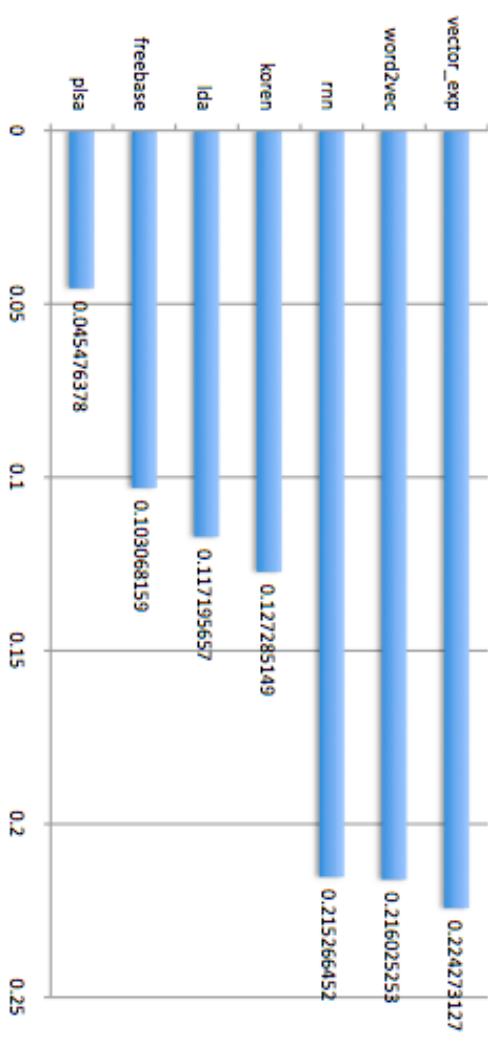


## deepwalk

# Playlists at Spotify

- ‘words’ are songs
- ‘sentences’ are playlists

Great performance on ‘related artists’



## Playlists at Spotify

# Fixes at Stitch Fix

Let's try:

‘words’ are styles  
‘sentences’ are fixes

# Fixes at Stitch Fix

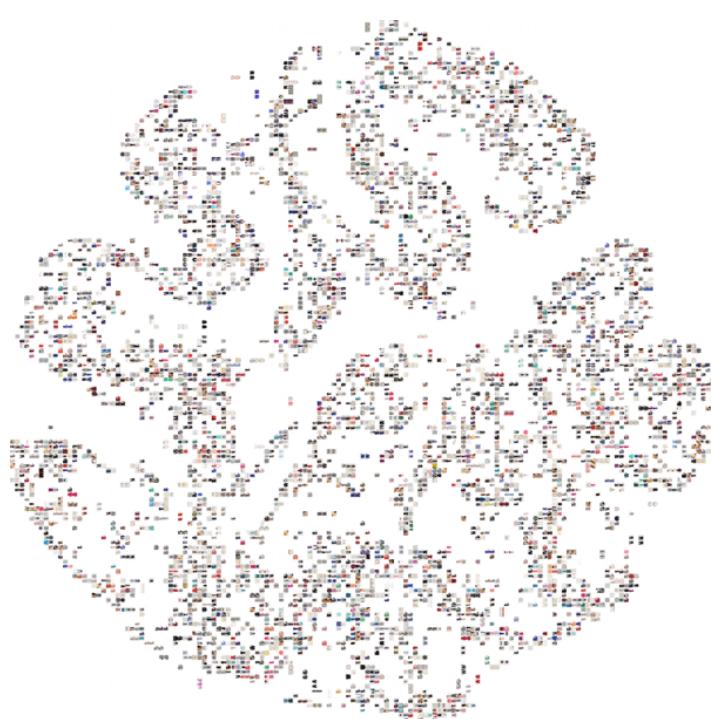
Learn similarity between styles

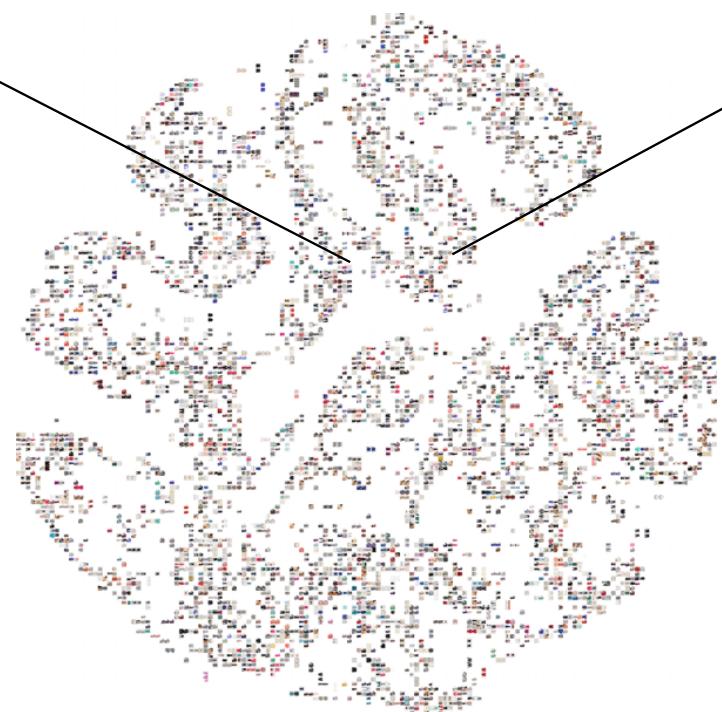
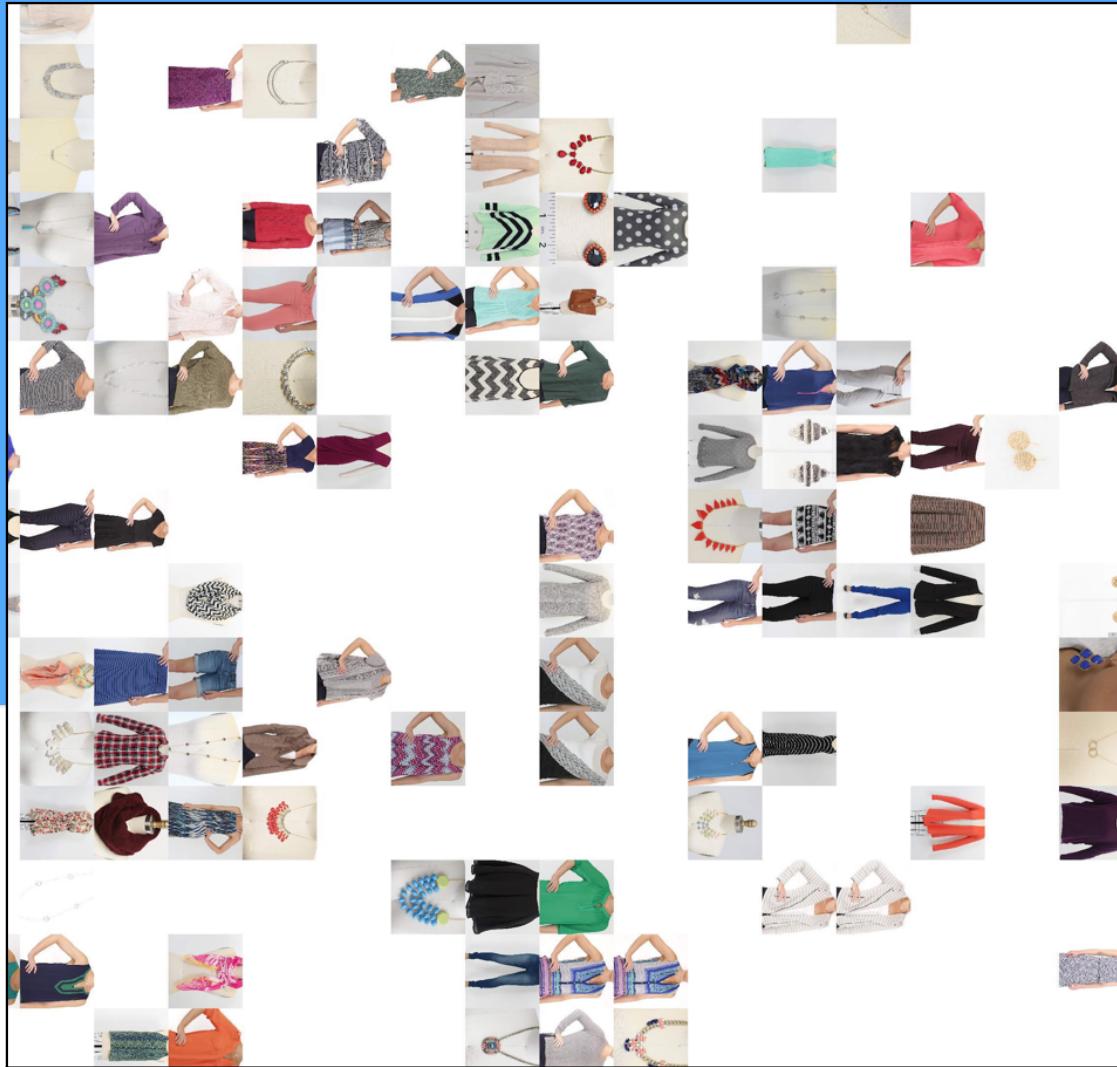
because they co-occur

Learn ‘coherent’ styles

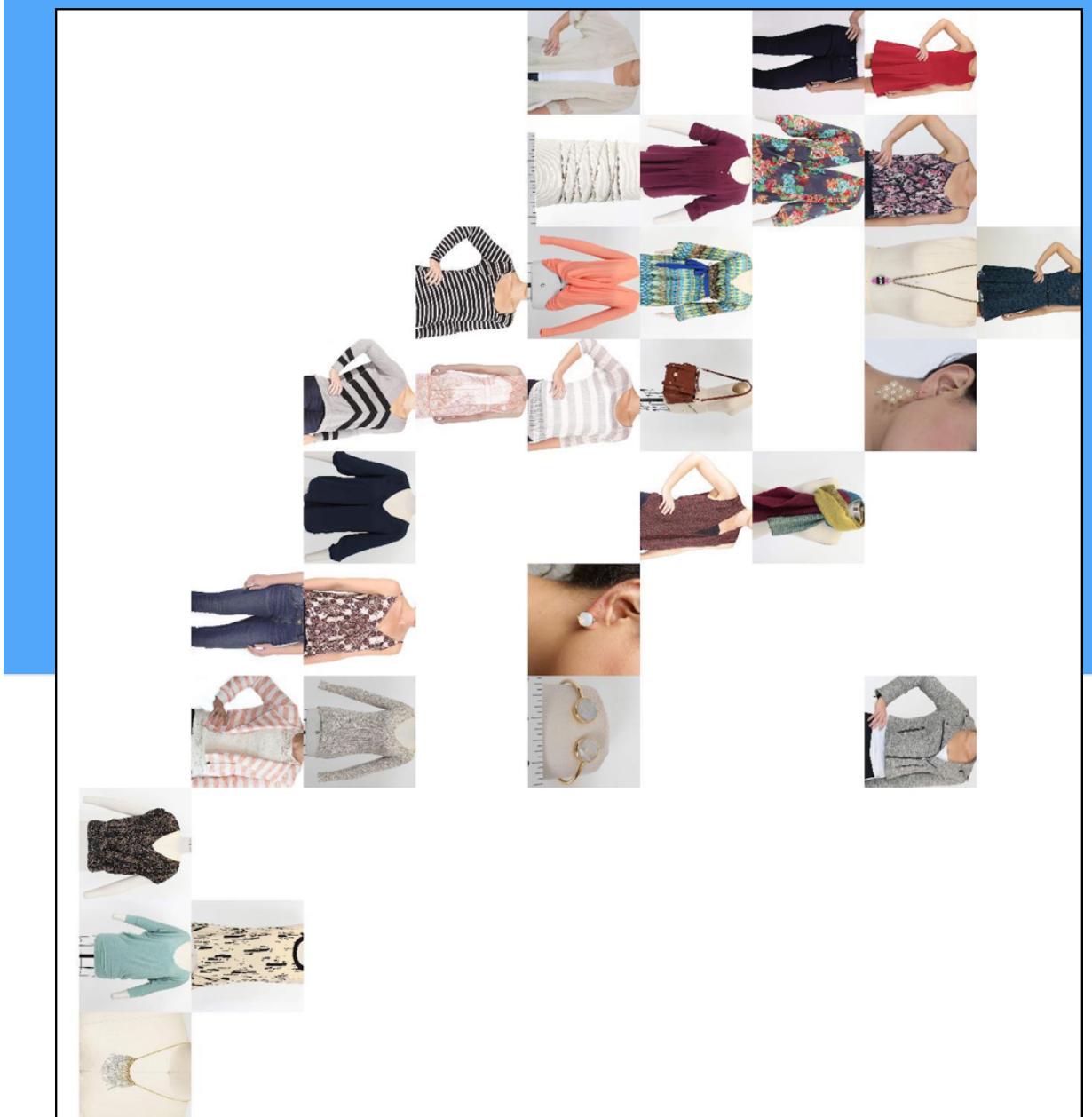
# Fixes at Stitch Fix?

Got lots of structure!

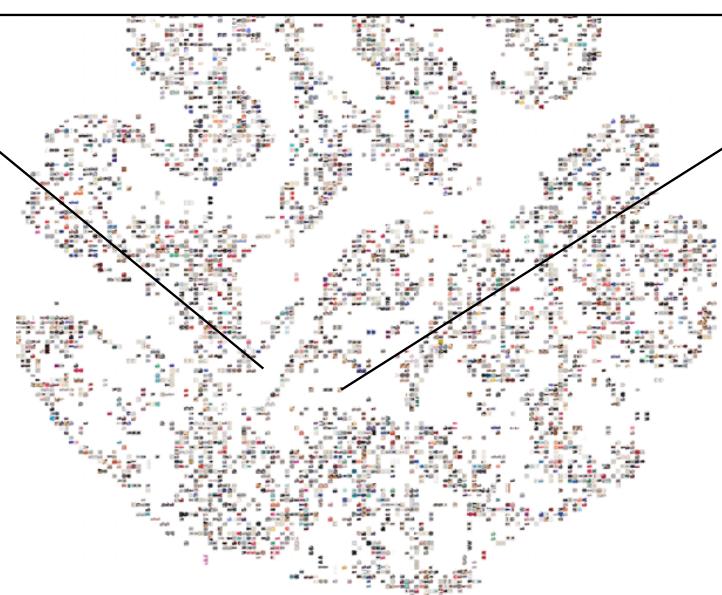




sequence  
learning



Nearby regions are  
consistent ‘closets’



sequence  
learning





A specific Ida2vec model

Our text blob is a comment that comes from a region\_id and a style\_id





$$L = \sigma(c * w) + \sigma(-c * w_{neg})$$

$$context = \vec{c_{ij}} = \vec{region_i} + \vec{style_j}$$

$$\vec{region_i} = \Sigma_{k=0}^{n\_topics} u_{ik} \cdot \vec{m_k}$$

$$\vec{style_j} = \Sigma_{l=0}^{n\_topics} u_{jl} \cdot \vec{n_l}$$

$$\vec{u}\sim dirichlet(\alpha_1)$$

$$\vec{v}\sim dirichlet(\alpha_2)$$

$$take\_rate\_in\_region \sim 5.0*\sigma(W\cdot \vec{u})$$

## The full likelihood model

$$L = \sigma(c * w) + \sigma(-c * w_{neg})$$

$$\text{context} = \vec{c}_{ij} = \vec{\text{region}}_i + \vec{\text{style}}_j$$

$$\vec{\text{region}}_i = \sum_{k=0}^{n\_topics} u_{ik} \cdot \vec{m_k}$$

$$\vec{\text{style}}_j = \sum_{l=0}^{n\_topics} u_{jl} \cdot \vec{n_l}$$

$$\vec{u} \sim \text{dirichlet}(\alpha_1)$$

$$\vec{v} \sim \text{dirichlet}(\alpha_2)$$

$$\text{take\_rate\_in\_region} \sim 5.0 * \sigma(W \cdot \vec{u})$$

$$L = \sigma(c * w) + \sigma(-c * w_{\text{neg}})$$

First part of the loss function is given **context** predict **word**.

**Don't** predict a **negative word**. These are words that are in our vocabulary somewhere, but not in our example.

We get negative samples **not** uniformly, but proportional to the word frequency<sup>3/4</sup> (yes, the <sup>3/4</sup> power is weird and ad hoc but totally works awesomely for word2vec)

$$L = \sigma(c * w) + \sigma(-c * w_{neg})$$

$$context = \vec{c}_{ij} = \vec{region}_i + style_j$$

Context is made up from more than one part -- many 'contexts' available.

In this case, instead of one document, we can have many regions, or styles.

In LDA, this context is a single term: the latent document vector that 'generates' words.

In word2vec, this context is the 'pivot' word. Word2vec picks a random 'context' word in the corpus, centers a window around it, and tries to predict other words within that context.

In both word2vec and LDA context is one term, either a document or a word. For lda2vec, we can have more than one term, we can have as many contexts as we like!

$$L = \sigma(c * w) + \sigma(-c * w_{neg})$$

$$context = \vec{c}_{ij} = \vec{region}_i + style_j$$

$$\vec{region}_i = \sum_{k=0}^{n\_topics} u_{ik} \cdot \vec{m_k}$$

$$\vec{style}_j = \sum_{l=0}^{n\_topics} u_{jl} \cdot \vec{n_l}$$

Each context (e.g., **region** or **style**) is decomposed into **topics vectors** and **weights** on those common **topics vectors**. One context has one shared set of topic vectors (think of these as cluster centroids) and every ‘document’ in that context (think of 1 of 50 states, 1 of 20k styles) has a weight/membership onto each of those topic vectors (think topics like northeast, midwest for region or tops, bottoms, boho, romantic for style topics)

This forces the context vectors onto **a limited set of basis vectors**. Interpret this set, and you can generalize what each region vector and style vector means. For example, one **topics vector** might be close to the **word vector** for ‘hand\_bag’, ‘purse’, ‘bag’ indicating that that topic is a handbags topic. And then anything with big **weight** in that topic might be a handbag.

$$L = \sigma(c * w) + \sigma(-c * w_{neg})$$

$$\text{context} = \vec{c_{ij}} = \vec{\text{region}_i} + \text{style}_j$$

$$\vec{\text{region}_i} = \Sigma_{k=0}^{n\_topics} u_{ik} \cdot \vec{n_k}$$

$$\vec{\text{style}_j} = \Sigma_{l=0}^{n\_topics} u_{jl} \cdot \vec{n_l}$$

$$\vec{u} \sim \text{dirichlet}(\alpha_1)$$

$$\vec{v} \sim \text{dirichlet}(\alpha_2)$$

But the weights can still end up being very dense -- which meant everyone of my documents was a mixture of almost every component. This made it difficult to interpret what the document was, because it had membership in many groups.

So next we enforce a simplex with dirichlet & enforce sparsity with the concentration on the **weights**. The dirichlet is also nice but not critical, we could've had a non-negative decomposition or just stuck with all reals. But since Dirichlet components sum to 100%, it is easier to explain to analysts that a document is “10% of some\_topic + 90% some\_other\_topic” rather than saying “-2.3 \* some\_topic and +0.5 of some\_other\_topic”.

$$L = \sigma(c * w) + \sigma(-c * w_{neg})$$

$$\text{context} = c_{ij}^{\vec{r}} = \vec{region}_i + \vec{style}_j$$

$$\vec{region}_i = \sum_{k=0}^{n_{topics}} u_{ik} \cdot \vec{n_k}$$

$$\vec{style}_j = \sum_{l=0}^{n_{topics}} u_{jl} \cdot \vec{n_l}$$

$$\vec{u} \sim dirichlet(\alpha_1)$$

$$\vec{v} \sim dirichlet(\alpha_2)$$

$$\text{take\_rate\_in\_region} \sim 5.0 * \sigma(W \cdot \vec{u})$$

Finally, we can make this ‘supervised’ by saying that the topic **weights** correlate through (matrix **W**) with some **target outcome**.



Let's make  $v_{DOC}$  into a mixture...

topic 1 = “religion”

*Milosevic*

*absentee*

*Indonesia*

*Lebanese*

*Isrealis*

*Karadzic*

$$v_{DOC} = 10\% \text{ religion} + 89\% \text{ politics} + \dots$$



topic 2 = “politics”

*Trinitarian*  
*baptismal*

*Pentecostals*

*bede*

*schismatics*

*excommunication*

