

Twitter User Classification



Devashish Deshpande

Undergraduate Student

RaRe Technologies Incubator Program



GitHub: dsquareindia

Blogs: <https://rare-technologies.com/blog/>

Gensim: Topic Modeling in python

RaRe-Technologies / **gensim** Unwatch 252 Unstar 2,978 Fork 1,087

[Code](#) [Issues 137](#) [Pull requests 39](#) [Projects 0](#) [Wiki](#) [Pulse](#) [Graphs](#)

Topic Modelling for Humans <http://radimrehurek.com/gensim/>

[2,595 commits](#) [11 branches](#) [35 releases](#) [123 contributors](#) [LGPL-2.1](#)

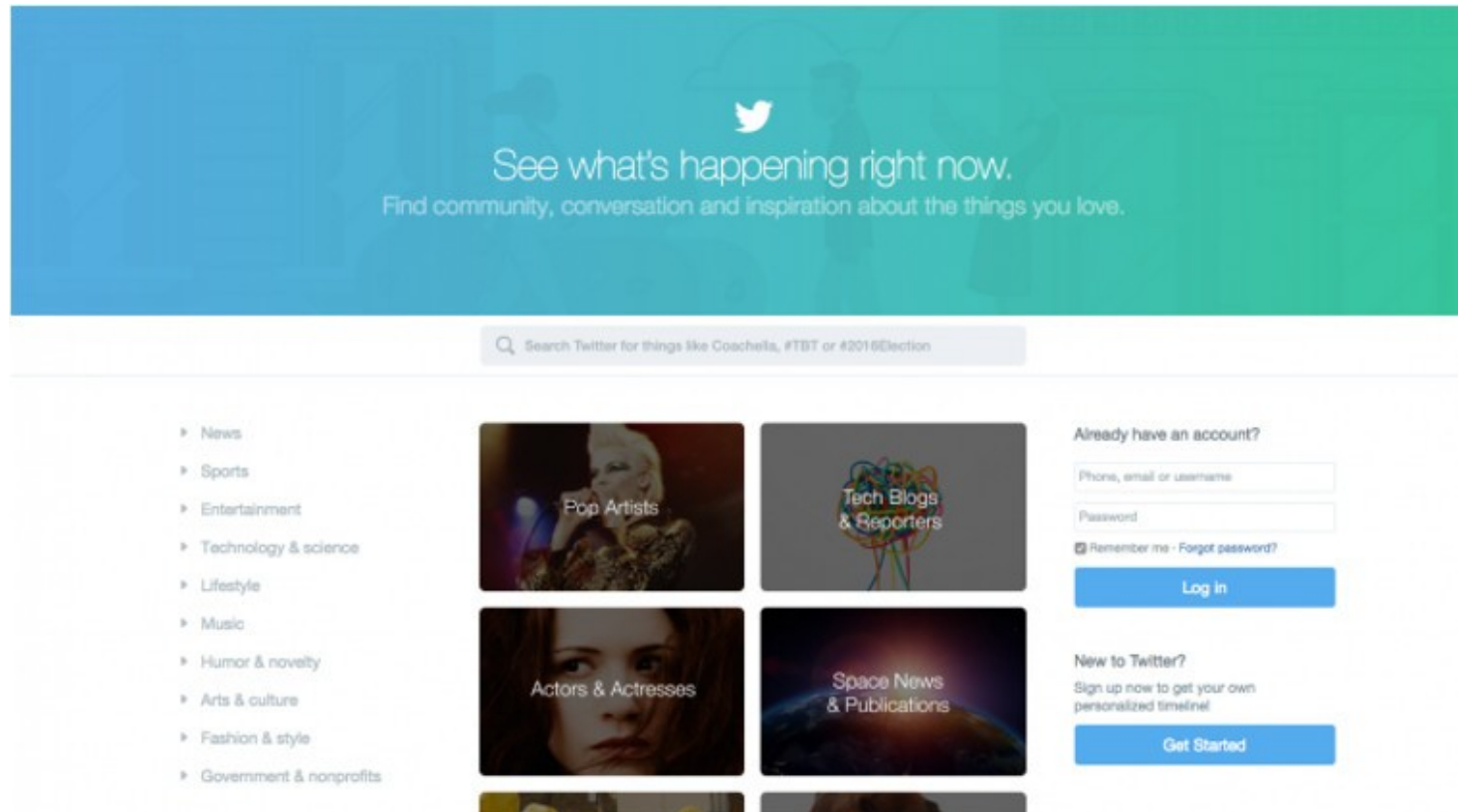
Branch: **develop** [New pull request](#) [Create new file](#) [Upload files](#) [Find file](#) [Clone or download](#)

tmylk committed on **GitHub** Merge pull request #864 from dsquareindia/supported_models Latest commit edc83c2 5 hours ago

continuous_integration/appveyor	find_links on rackspace	11 months ago
docs	Merge pull request #856 from einon/develop	6 hours ago
gensim	CoherenceModel documentation addition	2 days ago
.gitignore	add ipynb checkpoints into gitignore (#724)	4 months ago
.travis.yml	Annoy as an external similarity index for word2vec and doc2vec (#774)	3 months ago
CHANGELOG.md	Update CHANGELOG.md	28 days ago
CONTRIBUTING.md	project doc updates (CONTRIBUTING & README)	7 months ago
COPYING	LGPL v2.1 clarification	6 months ago
MANIFEST.in	Gensim adopters table added (#754)	3 months ago
README.md	Update README.md	2 months ago



Business Problem



- Millions of users on twitter
- Need to make category-wise recommendations
- How classify a user into a category?

Can you classify these users?

Soirée parfaite !   

50 goles en 81 partidos de Liga con el Atleti. El trabajo da sus frutos 

Nouvelles couleurs et toujours le même objectif : GOLAZO
 #GriziPums

Can you classify these users?

Soirée parfaite ! [F] [R] 🏆

50 goles en 81 partidos de Liga con el Atleti. El trabajo da sus frutos 💪

Nouvelles couleurs et toujours le même objectif : GOLAZO ⚽ #GriziPums



Can you classify these users?

The band have FOUR @mtvema nominations! Voting only takes a couple of clicks (and you can do as often as you like)

In case you missed it - 18 new #AHFODtour shows confirmed for US & Canada in Aug/Sep/Oct 2017

Face à la demande exceptionnelle @Coldplay ajoute une date supplémentaire le 16 juillet 2017 au @StadeFrance

Can you classify these users?

The band have FOUR @mtvema nominations! Voting only takes a couple of clicks (and you can do as often as you like)

In case you missed it - 18 new #AHFODtour shows confirmed for US & Canada in Aug/Sep/Oct 2017

Face à la demande exceptionnelle @Coldplay ajoute une date supplémentaire le 16 juillet 2017 au @StadeFrance



Can you classify these users?

Vous venez à la #PyConFr ? Soyez sympa: pensez à vous inscrire qu'on sache combien on sera: <http://2016.pycon.fr/>
!#Python #PyCon #Rennes

Regardez ces beaux t-shirts ! <https://2016.pycon.fr/t-shirts.html> A sérigraphier soi-même pour ceux et celles qui veulent, le jour de l'évènement !

Qui est de PyCon la semaine prochaine à Rennes? Je serais sur place jeudi. Un AFPyro quelque part le jeudi soir? #python #pyconfr

Can you classify these users?

Vous venez à la #PyConFr ? Soyez sympa: pensez à vous inscrire qu'on sache combien on sera:
<http://2016.pycon.fr/> !#Python
#PyCon #Rennes

Regardez ces beaux t-shirts !
<https://2016.pycon.fr/t-shirts.html> A sérigraphier soi-même pour ceux et celles qui veulent, le jour de l'évènement !

Qui est de PyCon la semaine prochaine à Rennes? Je serais sur place jeudi. Un AFPyro quelque part le jeudi soir? #python #pyconfr



Human mind is a great classifier!

How can we model our thought process?

Creating our dataset

- Created a toy dataset consisting of 8 categories
- Each category consists of:
 - Popular individuals in that category
 - Popular magazine/tv channel handles
- Both tweet in a different style
- Need to capture all kinds of tweets

Creating our dataset

Neil deGrasse Tyson

Not that anybody asked, but on Mercury, I'm 240 years old, and on Saturn, I'm just 2.

If ComicCon people ruled the world, international conflicts would be resolved entirely by plastic light saber fights in bars

NASA

What science experiments did crew members on @Space_Station work on this week? Find out here

We're reducing carbon emissions from aviation by creating aircraft designed to dramatically reduce fuel use & noise: <http://go.nasa.gov/2e1BRpn>

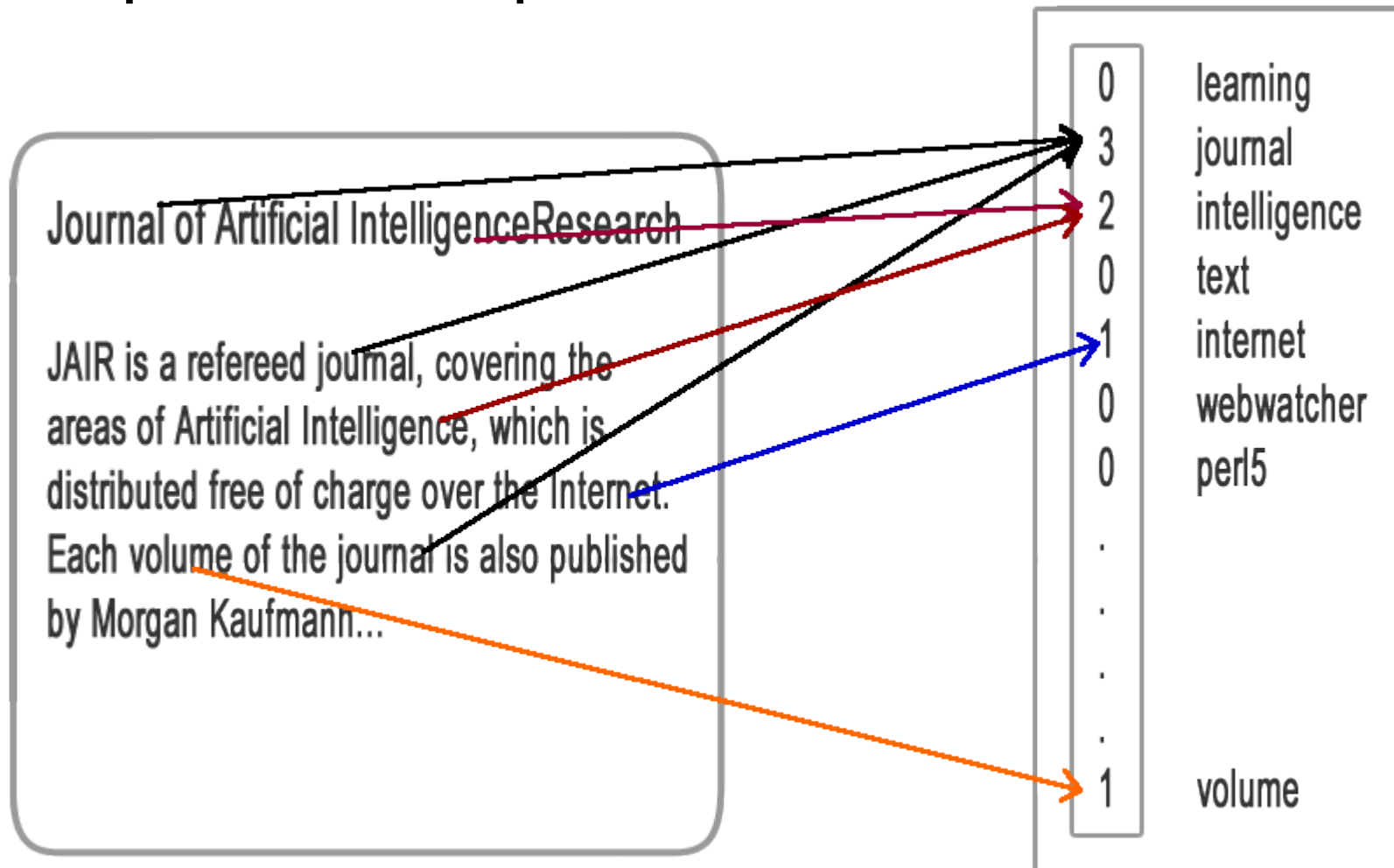
Both should be classified into “Science”

Creating our dataset

- Our categories:
 - Tech
 - Business & CEOs
 - Entertainment
 - Science
 - Fashion, Travel and Lifestyle
 - Sports
 - Music
 - Politics
- Music
 - Iron Maiden
 - David Guetta
 - Rolling Stone magazine
 - MTV Music and more...

Bag of words

- Simple count representation



Bag of words

- Can be great in our case
- Tweets have distinct words which are used often

```
In [49]: most_influential_words(clf_count, count_vectorizer, category_index=7) # Top words for science category
```

```
Out[49]: [u'science',  
          u'earth',  
          u'scientist',  
          u'space',  
          u'alien',  
          u'study',  
          u'dystrophaeus_fossiltime',  
          u'bone',  
          u'book',  
          u'astronaut']
```

Image from accompanying notebook

TF-IDF

- Slightly more complex
- Normalizing using number of docs the term appears in

$$w_{x,y} = tf_{x,y} \times \log \left(\frac{N}{df_x} \right)$$

TF-IDF

Term x within document y

$tf_{x,y}$ = frequency of x in y

df_x = number of documents containing x

N = total number of documents

Can we try something else?

We are trying to classify documents into topics.
Can we find out the topics using topic modeling
and then classify?

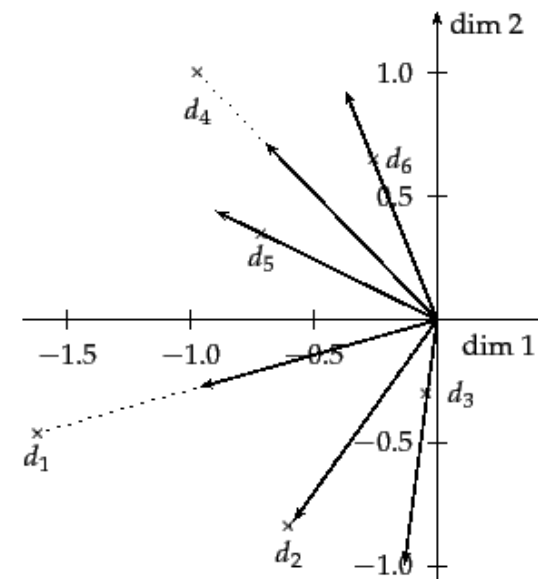
Topic modeling

- Gensim has a number of topic modeling techniques available! We'll be using
 - Latent Semantic Indexing
 - Hierarchical Dirichlet Process
 - Latent Dirichlet Allocation

Latent Semantic Indexing

- Uses SVD
- Maps terms and documents to a latent semantic space

		d_1	d_2	d_3	d_4	d_5	d_6	
	ship	1	0	1	0	0	0	
	boat	0	1	0	0	0	0	
	ocean	1	1	0	0	0	0	
	voyage	1	0	0	1	1	0	
	trip	0	0	0	1	0	1	



Latent Semantic Indexing

- Can rank topics automatically
- Need to provide num_topics: the number of latent dimensions

```
In [8]: lsi.print_topics(2)
```

```
Out[8]: [(0,  
          u'0.703*"trees" + 0.538*"graph" + 0.402*"minors" + 0.187*"survey" + 0.061*"system" + 0.060*"ti  
me" + 0.060*"response" + 0.058*"user" + 0.049*"computer" + 0.035*"interface"'),  
         (1,  
          u'-0.460*"system" + -0.373*"user" + -0.332*"eps" + -0.328*"interface" + -0.320*"response" + -  
0.320*"time" + -0.293*"computer" + -0.280*"human" + -0.171*"survey" + 0.161*"trees"')]
```

Image taken from [Topics_and_Transformations notebook](#) trained on deerwester.mm corpus

Hierarchical Dirichlet Process

- Fully unsupervised!
- No need of num_topics parameter
- Determines number of topics through posterior inference
- Non-parametric generalization of LDA
- Can sometimes fail to capture granularity in topics

Latent Dirichlet Allocation

- “Celebrity topic model”. One of the most popular topic modeling algorithms
- Generative process
- Each document is a mixture of topics
- Each topic generates words
- Needs num_topics for fitting

Topic-Word coloring with LDA

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which

What is a good LDA model?

- Come up with interpretable topics
- Infer topic distribution

Football LDA model

Topic 0: mourinho, red_devils, old_trafford, bad_team...

Topic 1: wenger, henry, invincibles.....

Topic 2: aguero, etihad, england, premier_league

Topic 3: blues, football, roman, bridge

We have so many topic models. Great!

But how do we compare them?

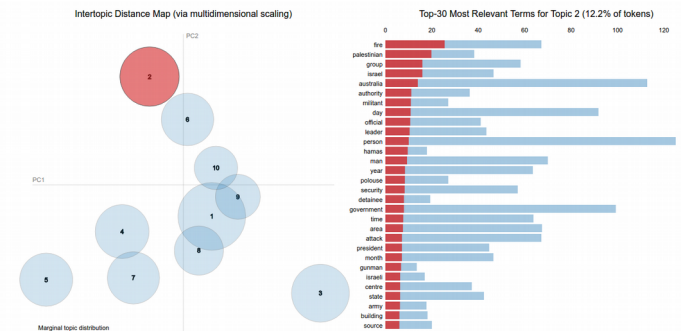
[illegible]

- Qualitative analysis
- Can get ugly

```

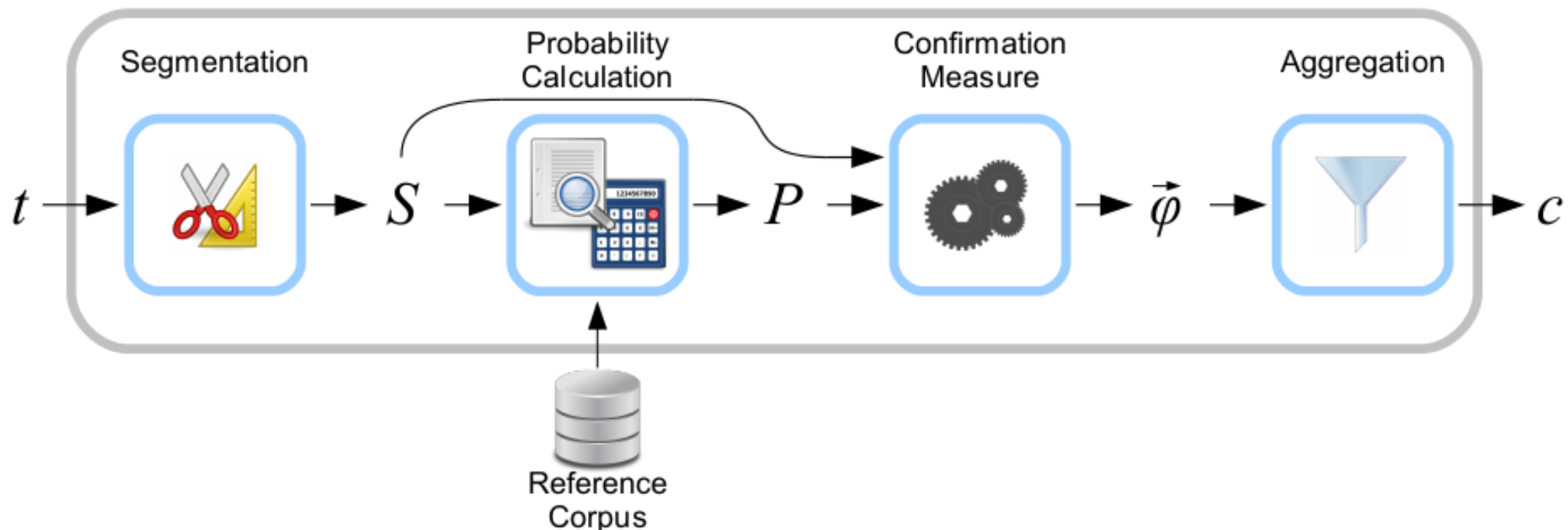
u'topic 8: 0.003*india + 0.003*album + 0.003*cheese + 0.002*chocolate + 0.002*cake + 0.002*summit + 0.002*watch +
0.002*song + 0.002*person + 0.002*ice_cream + 0.002*minister + 0.002*music + 0.002*water + 0.002*talk + 0.001*tick
et + 0.001*david_bowie + 0.001*pizza + 0.001*time + 0.001*fan + 0.001*president',
u'topic 9: 0.016*travel + 0.013*recipe + 0.006*ironmaiden + 0.003*rio + 0.002*thank + 0.002*time + 0.002*india +
0.002*ironmaiden_spain + 0.002*woman + 0.002*year + 0.002*day + 0.002*tonight + 0.002*ironmaiden_italy + 0.002*tea
m + 0.002*cheese + 0.001*birthday + 0.001*china + 0.001*solo + 0.001*city + 0.001*congratulation',
u'topic 10: 0.004*watch + 0.004*goal + 0.003*theater + 0.003*india + 0.003*slvaus + 0.003*england + 0.002*pokemon
+ 0.002*australia + 0.002*sri_lanka + 0.002*game + 0.002*trainer + 0.002*team + 0.002*sky_sport + 0.002*man + 0.00
2*match + 0.002*mnf + 0.002*everton + 0.002*pakistan + 0.002*test + 0.002*pokemongo',
u'topic 11: 0.006*video + 0.004*one + 0.004*today + 0.004*thank + 0.003*apple + 0.003*day + 0.003*time + 0.003*wor
ld + 0.003*euro + 0.002*person + 0.002*bigtheparty + 0.002*party + 0.002*everyone + 0.002*night + 0.002*team + 0.00
1*coloring_book + 0.001*life + 0.001*euro_davidguetta + 0.001*friend + 0.001*tomorrow',
u'topic 12: 0.011*travel + 0.005*world + 0.004*cnnfood + 0.002*earth + 0.002*mission + 0.002*space + 0.002*way +
0.002*today + 0.002*city + 0.002*time + 0.002*spacecraft + 0.002*sample_return + 0.002*launch + 0.002*dessert + 0.
002*guide + 0.002*day + 0.002*japan + 0.001*view + 0.001*dish + 0.001*park',
u'topic 13: 0.004*album + 0.004*song + 0.004*video + 0.002*track + 0.002*share + 0.002*year + 0.002*watch + 0.002*
film + 0.001*review + 0.001*stream + 0.001*week + 0.001*science + 0.001*way + 0.001*life + 0.001*star_trek + 0.001*
recording + 0.001*person + 0.001*band + 0.001*book + 0.001*game',
u'topic 14: 0.004*goal + 0.003*time + 0.002*today + 0.002*day + 0.002*ageofultron_presstour + 0.002*chelsea + 0.00
1*game + 0.001*dystrophaeus_fossiltime + 0.001*diego_costa + 0.001*arsenal + 0.001*watch + 0.001*way + 0.001*man +
0.001*liverpool + 0.001*shirt + 0.001*year + 0.001*season + 0.001*fossilfriday + 0.001*presstour + 0.001*point',
u'topic 15: 0.003*space + 0.003*earth + 0.003*year + 0.002*recipe + 0.002*voyager + 0.002*scientist + 0.002*today
+ 0.002*travel + 0.002*day + 0.002*system + 0.002*pic + 0.002*time + 0.001*water + 0.001*star + 0.001*fossil + 0.0
01*summer + 0.001*edge + 0.001*tip + 0.001*idea + 0.001*tonight',
u'topic 16: 0.005*video + 0.005*canada + 0.004*album + 0.002*song + 0.002*singer + 0.002*china + 0.002*chine + 0.0
02*today + 0.002*frontman + 0.001*avec + 0.001*discussion + 0.001*shanghai + 0.001*pour + 0.001*metallica + 0.001*t
ime + 0.001*thank + 0.001*tour + 0.001*paralympique + 0.001*music + 0.001*day',
u'topic 17: 0.004*halloween + 0.003*year + 0.003*video + 0.003*vmas + 0.002*today + 0.002*time + 0.002*vma + 0.002
*fan + 0.002*game + 0.002*performance + 0.002*thank + 0.002*day + 0.002*ticket + 0.002*night + 0.002*tonight + 0.00
2*resetind + 0.002*halloween_destination + 0.002*kanye + 0.001*guy + 0.001*techno',
u'topic 18: 0.004*thank + 0.003*time + 0.003*show + 0.003*video + 0.003*tonight + 0.003*night + 0.002*year + 0.002

```



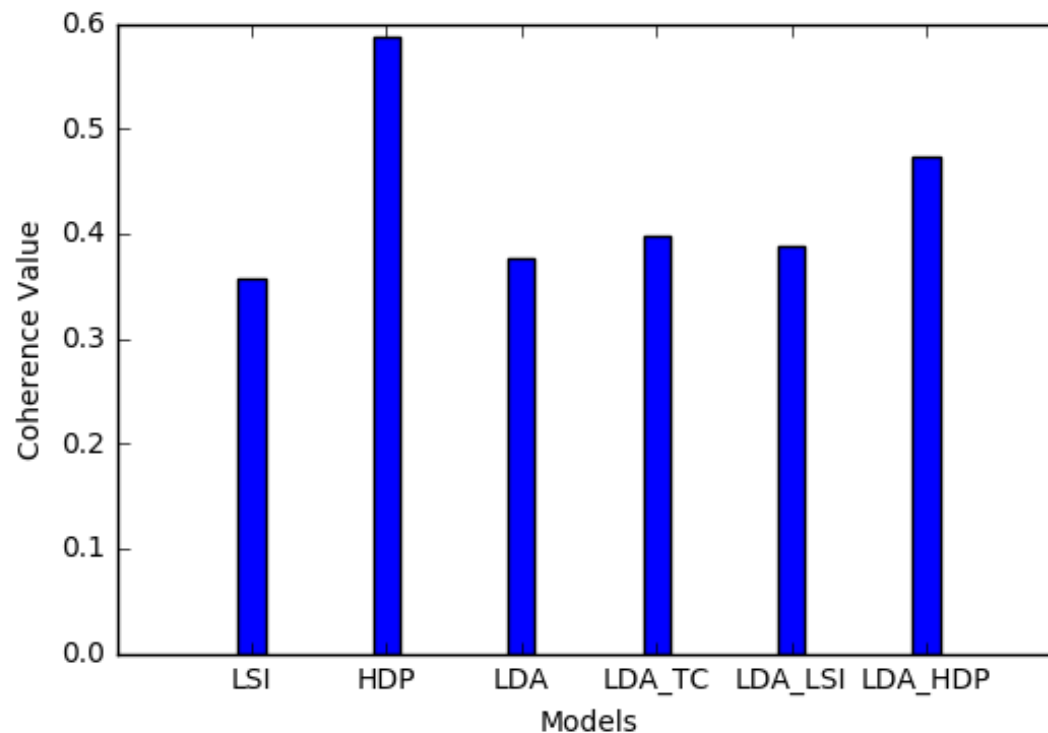
Topic Coherence

- Recently added to gensim
- Based on [this paper by Roeder et al](#)



Topic Coherence

- Quantitative
- Enter topics 't' -> Get coherence value 'c'
- Comparison much easier!



Topic Coherence

- Many other use-cases
 - Selecting num_topics for LDA
 - Rank topics within LDA model (pseudo-LSI)
 - Filtering LDA to improve model

Topic Coherence

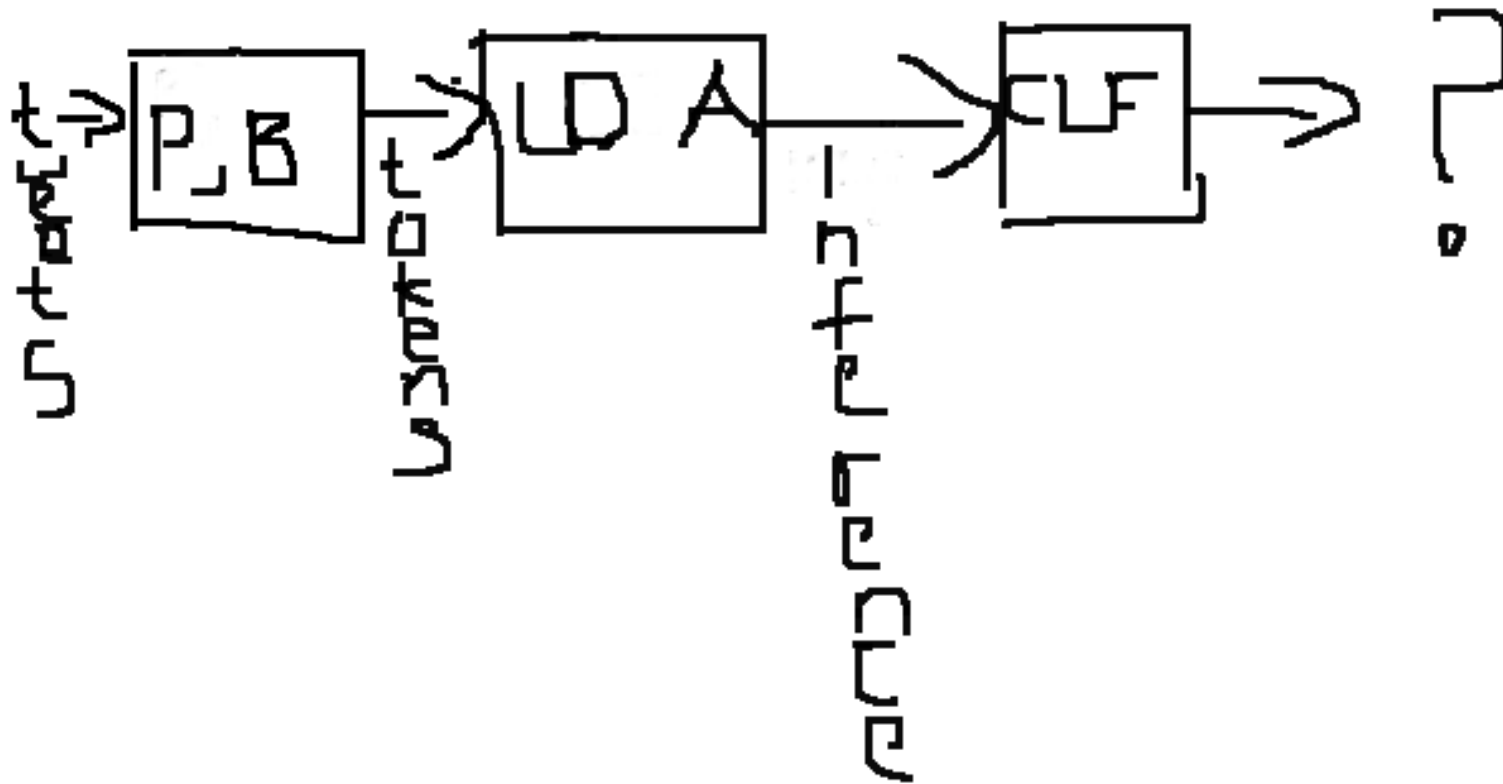
- Better LDA -> Better topics -> Better classification

```
[[('actor', 0.034688196735986693),  
  ('picture', 0.023163878883499418),  
  ('award', 0.023163878883499418),  
  ('comedy', 0.023163878883499418),  
  ('globe', 0.023163878883499418),  
  ('nomination', 0.023163878883499418),  
  ('actress', 0.023163878883499418),  
  ('film', 0.023163878883499418),  
  ('drama', 0.011639561031012149),  
  ('winner', 0.011639561031012149)],  
 [('virus', 0.064292949289013482),  
  ('user', 0.048074573973209883),  
  ('computer', 0.040350900997751814),  
  ('company', 0.028173623478117912),  
  ('email', 0.022580226976870982),  
  ('worm', 0.020928236506996975),  
  ('attachment', 0.014534311779706417),  
  ('outlook', 0.01260706654637953),  
  ('software', 0.011909411409069969),  
  ('list', 0.0088116041533348403)],  
  [('australia', 0.038230610979973961),  
   ('test', 0.03039802044037989),  
   ('day', 0.026478028361575149),  
   ('adam', 0.023237227270639361),  
   ('wicket', 0.018060239149805601),  
   ('match', 0.015652900511647725),  
   ('gilchrist', 0.015206348827236857),  
   ('steve_waugh', 0.01496754571623464),  
   ('south_africa', 0.013902623982144873),  
   ('selector', 0.012332915474867073)],
```

Top topics from [topic modeling tutorial on Lee corpus](#)

LDA for dimensionality reduction

- We'll be using our best LDA model for dimensionality reduction

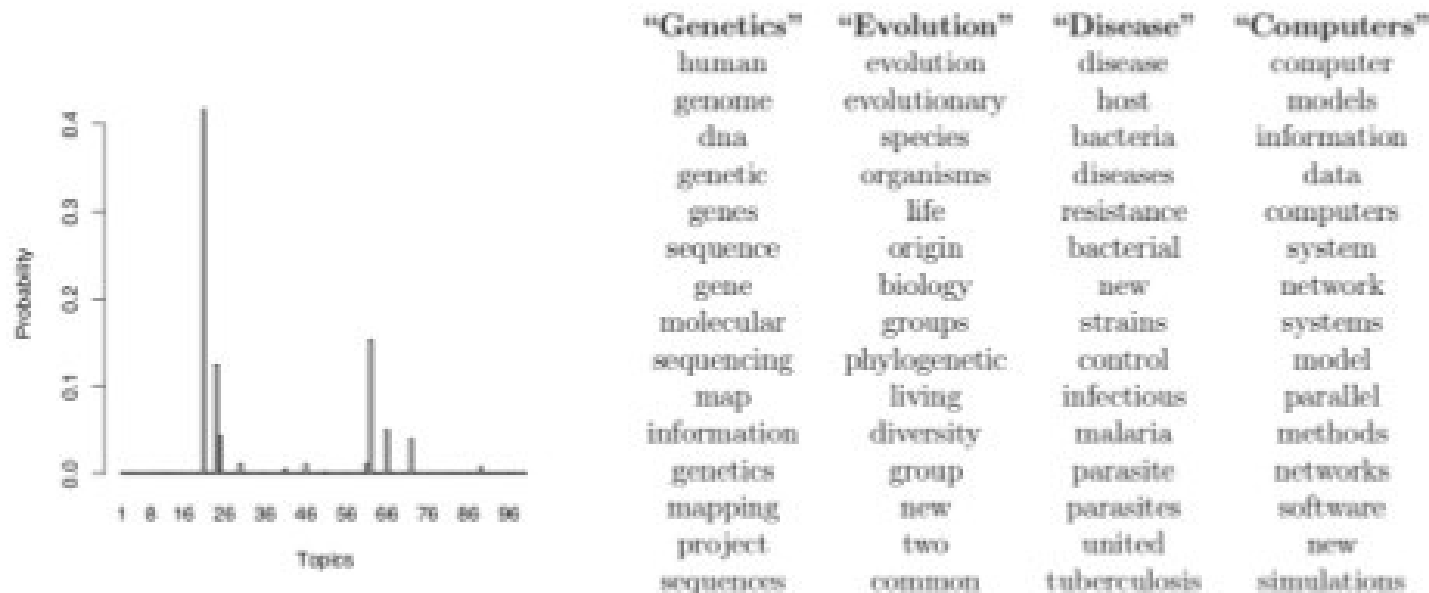


LDA for dimensionality reduction

- Using LDA inference to reduce dimensions to num_topics

Real inference with LDA

A 100-topic LDA model was fitted to **17,000 articles from the *Science* journal**.
At right are **the top 15 most frequent words** from the most frequent topics.
At left are the **inferred topic proportions** for the example article from previous slide.



GloVe Vectors

- Algorithm to map words to vectors!

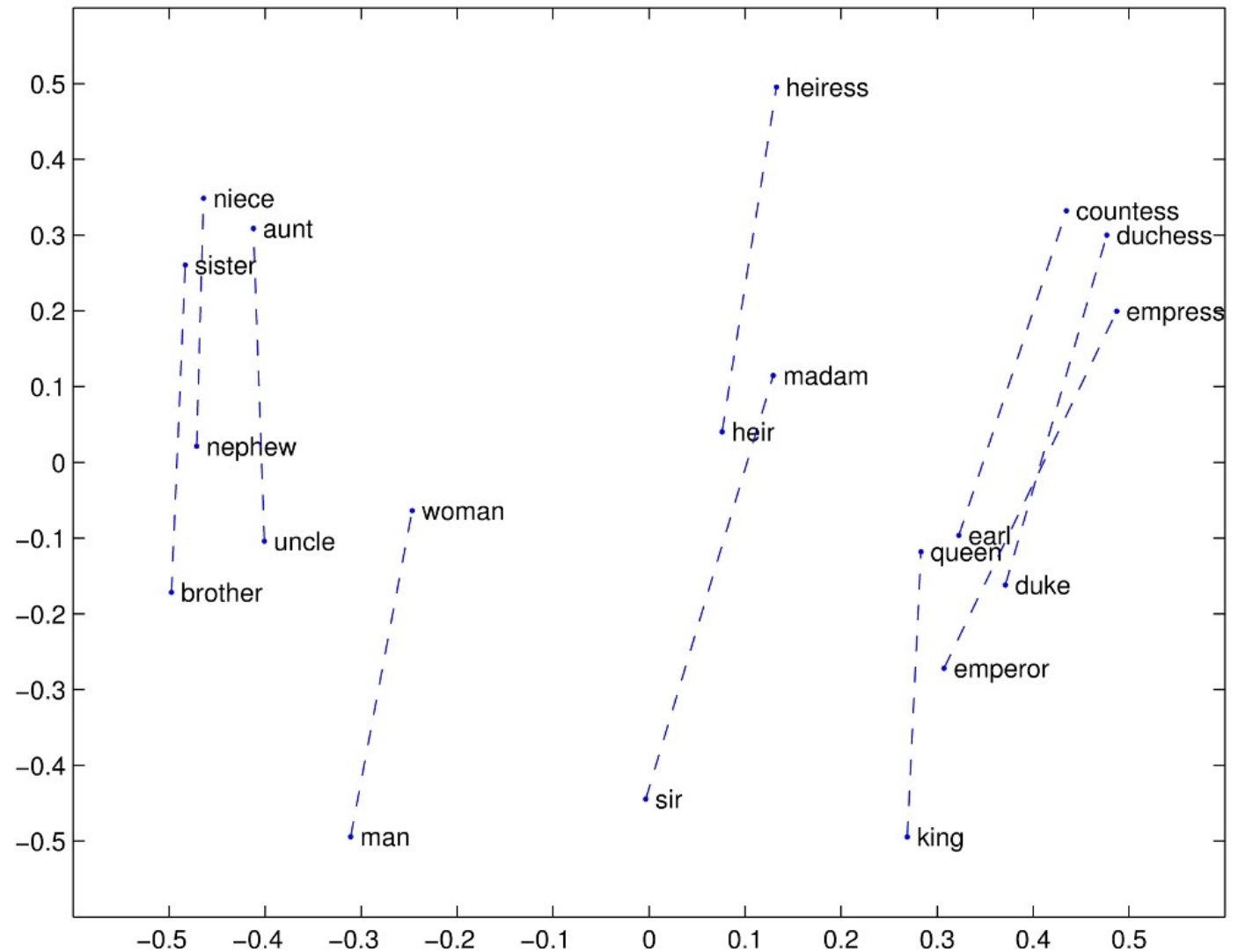
```
In [105]: ww.most_similar(positive=['arsenal'])
```

```
Out[105]: [(u'chelsea', 0.9786632061004639),  
            (u'liverpool', 0.9681325554847717),  
            (u'utd', 0.9317814707756042),  
            (u'united', 0.9310426712036133),  
            (u'manchester', 0.9226013422012329),  
            (u'tottenham', 0.9214313626289368),  
            (u'everton', 0.919910728931427),  
            (u'barca', 0.9044719934463501),  
            (u'swansea', 0.9019159078598022),  
            (u'qpr', 0.8937806487083435)]
```

Screenshot from accompanying notebook

GloVe Vectors

- Can do cool stuff!



Picture from [glove website](#)

Word2Vec: Makes your heart skip a gram

- How is GloVe different from Word2Vec?

GloVe

- Learns vectors from word-word co-occurrences
- “Reconstruction Loss”: Should keep variance while reducing dimensions

Word2Vec

- “Deep learning” predictive model
- CBOW, skip-gram
- Predicting target word given context

Word2Vec vs GloVe

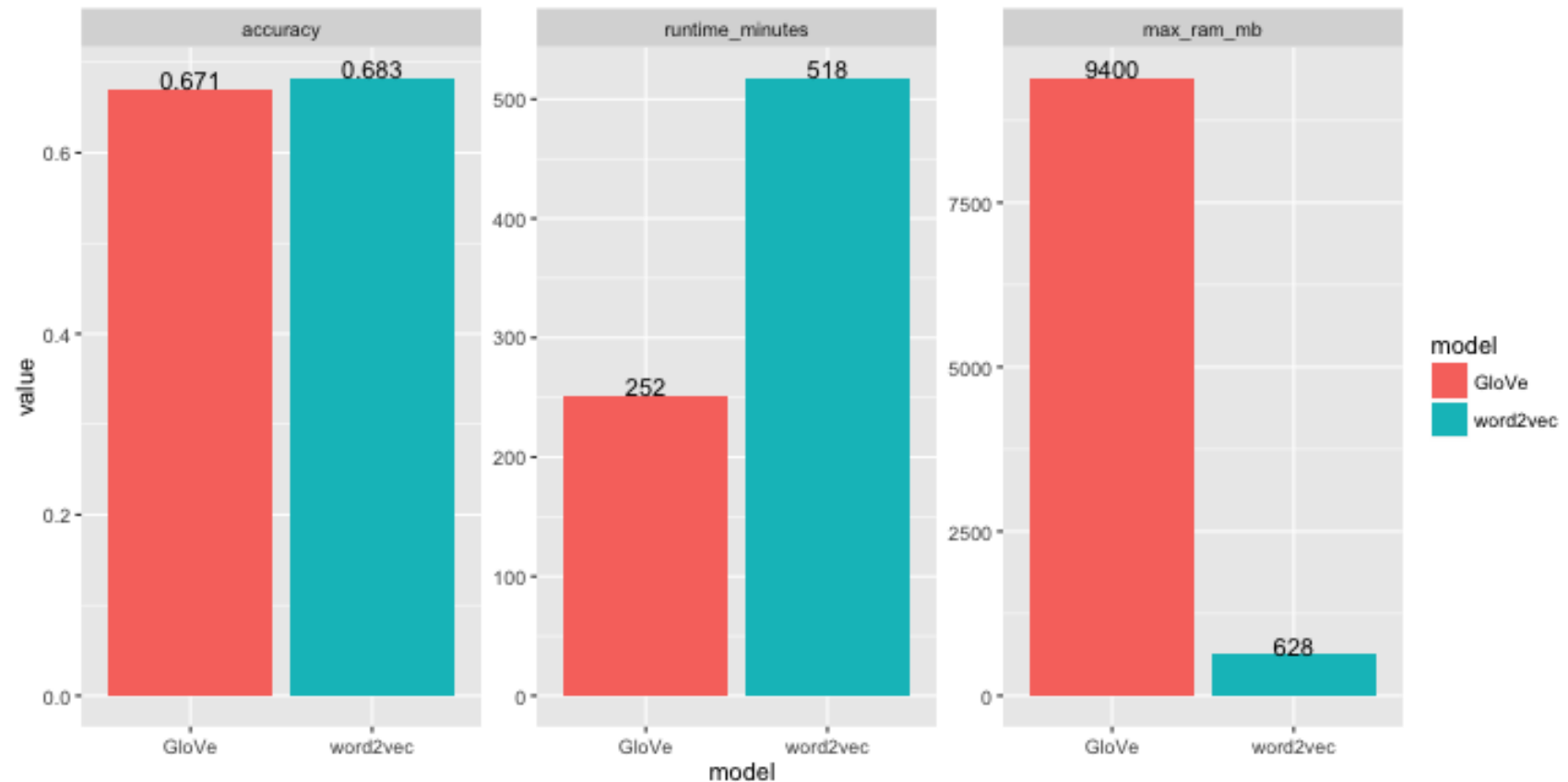


Image from <http://dsnotes.com/articles/glove-enwiki>

Using GloVe

- Our training time is zero!
- Using pre-trained twitter GloVe vectors by Stanford
- Converting to word2vec format for gensim compatibility using glove2word2vec script

Let's move on to the **tutorial!**