



Laboratoire de Sciences Cognitives et Psycholinguistique  
Ecole Normale Supérieure, Paris

# A Corpus-based Evaluation Method For Distributional Semantic Models

Abdellah Fourtassi & Emmanuel Dupoux

ACL- SRW, Sofia 2013



# Evaluation Methods for DSMs

---

## Extrinsic Methods

**Quantitatively** : Compare to human judgement (Benchmarks : Word association norms (Nelson et al, 1996), TOEFL synonym test (Laudauer, 1997)...) )

**Qualitatively** : Nearest neighbors ( researcher intuition )

HYPOTHESIS  
EXPERIMENT  
SCIENTIFIC  
OBSERVATIONS  
SCIENTISTS  
EXPERIMENTS  
SCIENTIST  
EXPERIMENTAL  
**TEST**  
METHOD

STUDY  
**TEST**  
STUDYING  
HOMEWORK  
NEED  
**CLASS**  
MATH  
TRY  
TEACHER  
WRITE

**CLASS**  
MARX  
ECONOMIC  
CAPITALISM  
CAPITALIST  
SOCIALIST  
SOCIETY  
SYSTEM  
**POWER**  
RULING

ENGINE  
FUEL  
ENGINES  
STEAM  
GASOLINE  
AIR  
**POWER**  
COMBUSTION  
DIESEL  
EXHAUST

Griffiths et al. (2007)

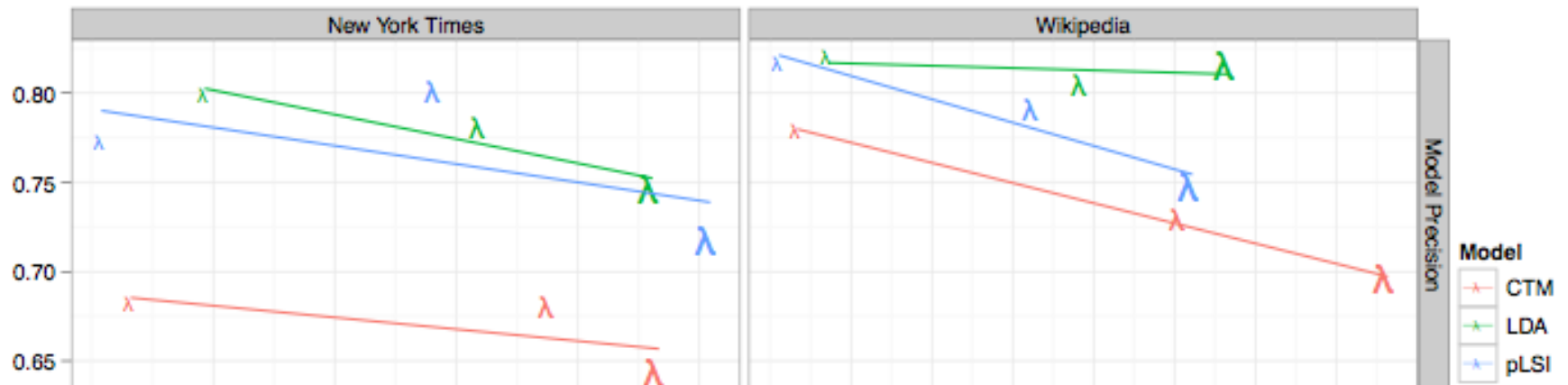
What if the researcher does not speak the language?

# Evaluation Methods for DSMs

## Intrinsic Methods

### Held-out Likelihood

- Possible only for probabilistic models (pLSA, LDA,..)
- Costly computationally
- Do not predict human judgement !! ( Chang, 2009)



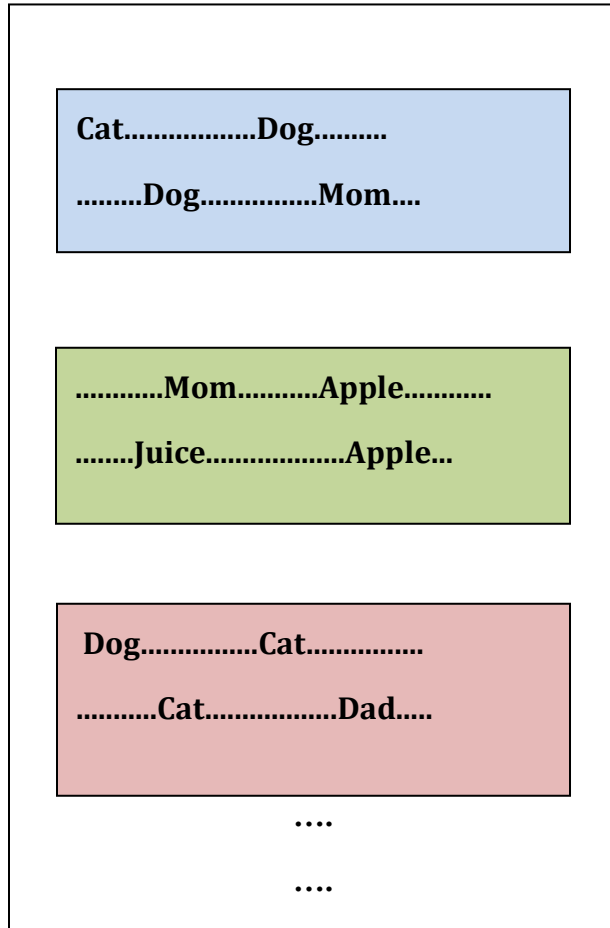
# Evaluation Method

---

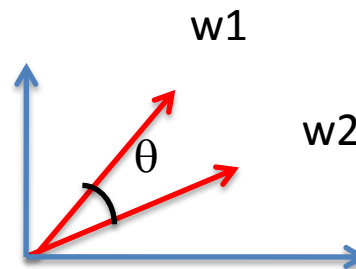
## Our Method

- Intrinsic : no external resources or mastery of the language required
- Predicts human judgement
- Easy to implement

# Latent Semantic Analysis



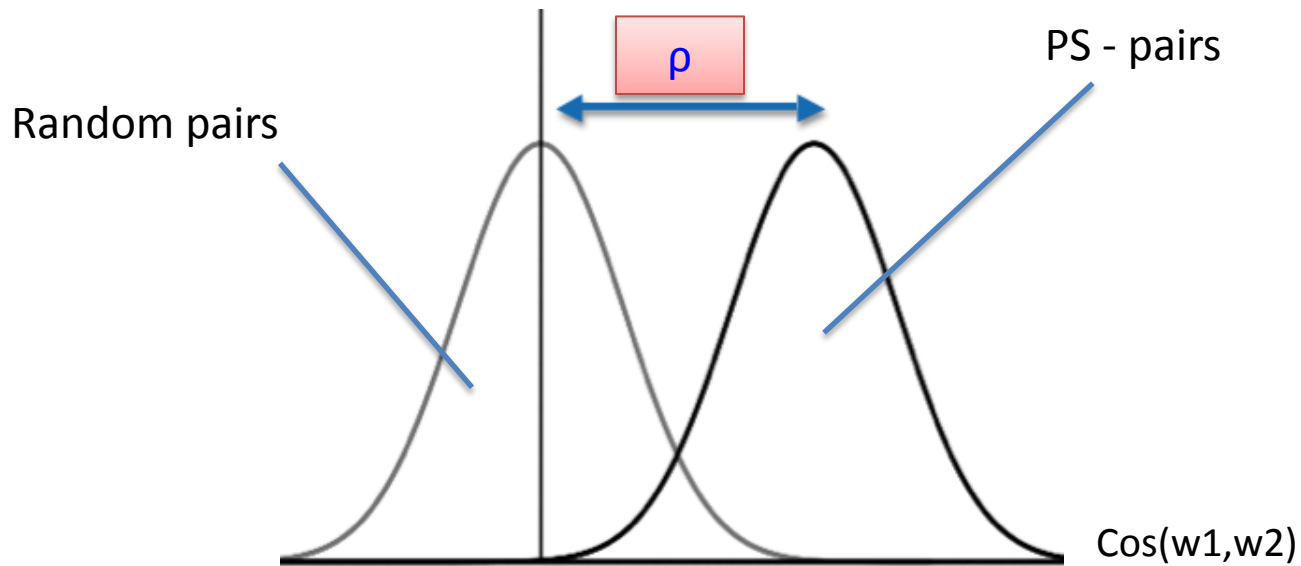
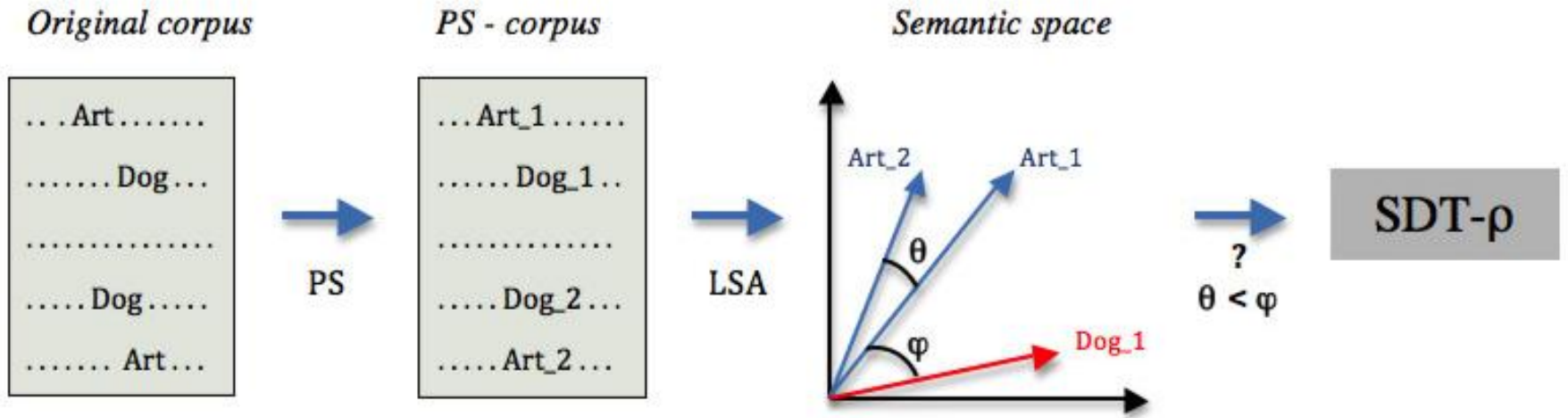
	<b>Context 1</b>	<b>Context 2</b>	<b>Context 3</b>	
Cat	<b>1</b>	<b>0</b>	<b>2</b>	...
Mom	<b>1</b>	<b>1</b>	<b>0</b>	...
Dog	<b>2</b>	<b>0</b>	<b>1</b>	...
Juice	<b>0</b>	<b>1</b>	<b>0</b>	...
Apple	<b>0</b>	<b>2</b>	<b>0</b>	...
Dad	<b>0</b>	<b>0</b>	<b>1</b>	...
...	...	...	...	...



$$\cos(w_1, w_2) = \frac{w_1^T w_2}{\|w_1\| \|w_2\|}$$

Landauer and Dumais (1996)

# SDT- $\rho$



# Evaluation Method

---

SDT- $\rho$



→ Intrinsic : no external resources or mastery of the language required

→ Predicts human judgement

→ Easy to implement



# Experiment I

---

## Word association Norms

	1	2	3
book	read	school	study

...

	1	2	3
Young	old	child	Happy

# Experiment I

## Word association Norms

book	1	2	3
	read	school	study
	...		
Young	1	2	3
	old	child	Happy
	...		

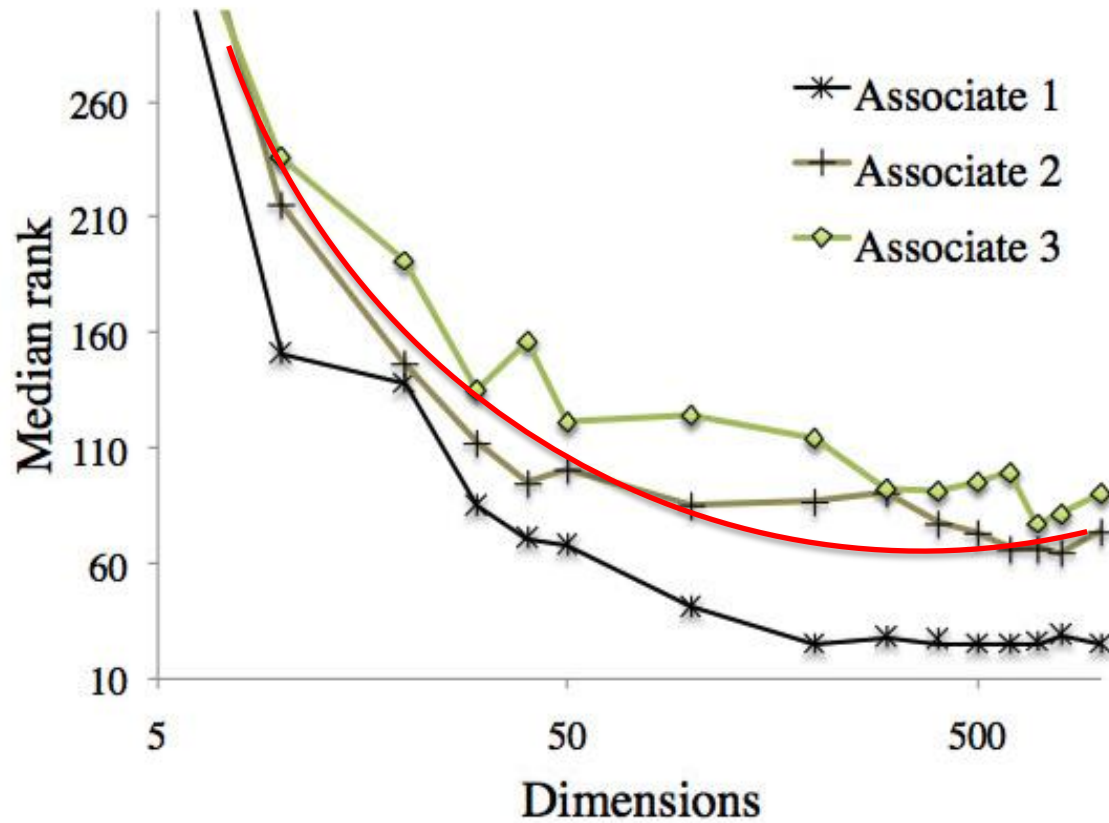
## LSA on Wikicorpus

Book	...	11	...	27	...	50
	...	read	...	school	...	study
	...	...				
Young	...	7	...	15	...	48
	...	old	...	child	...	happy
	...	...				

Median Rank  
of 1st  
associate

Median Rank  
of 1st  
associate

# Experiment I



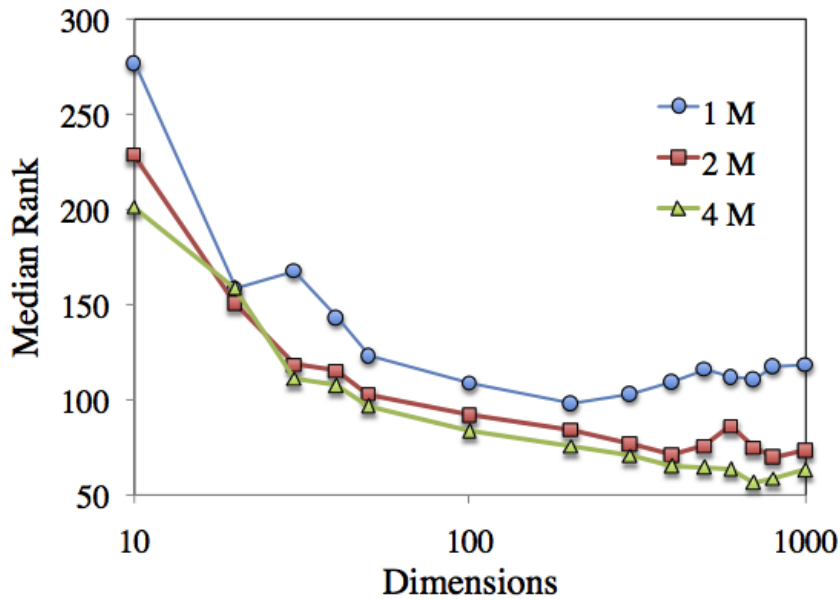
Median Rank



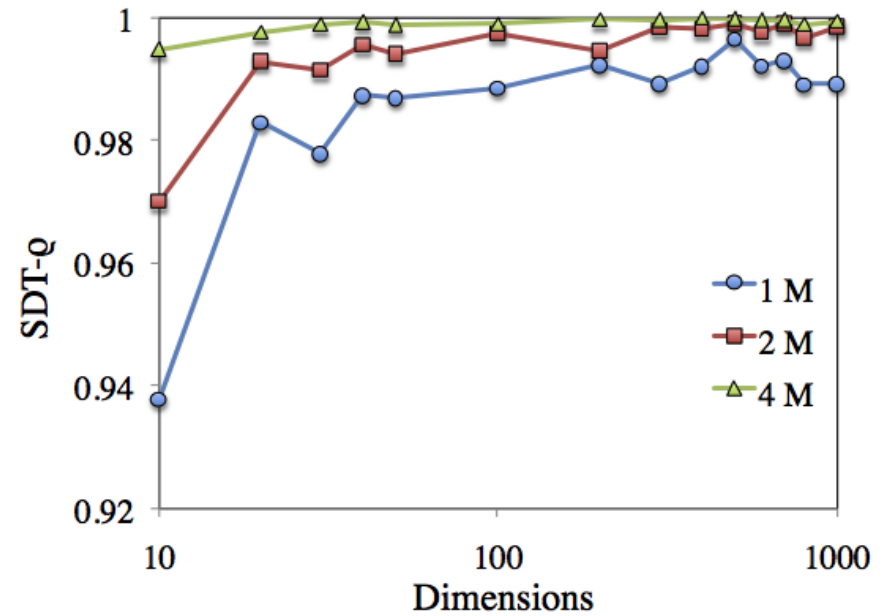
SDT- $\rho$

# Results I

Correlation :  
Dimensions and Corpus size



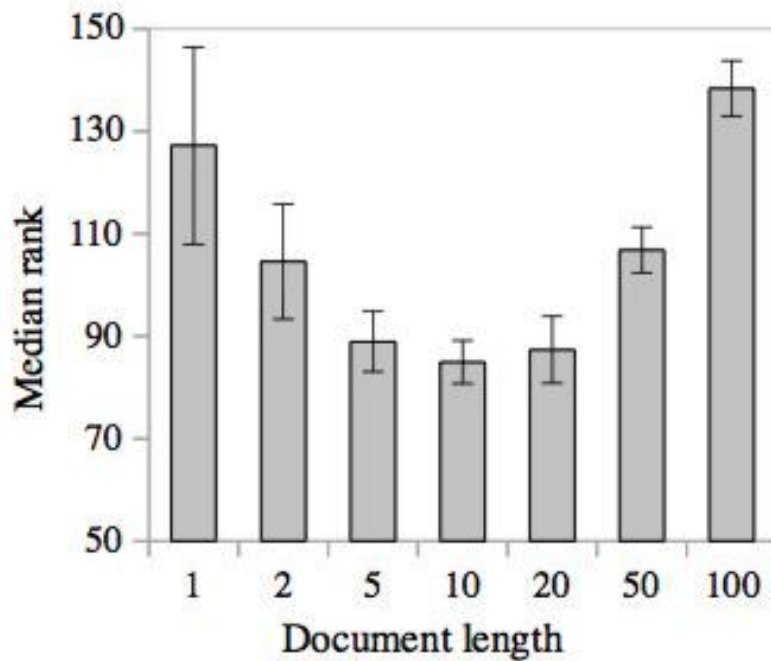
Lower is better



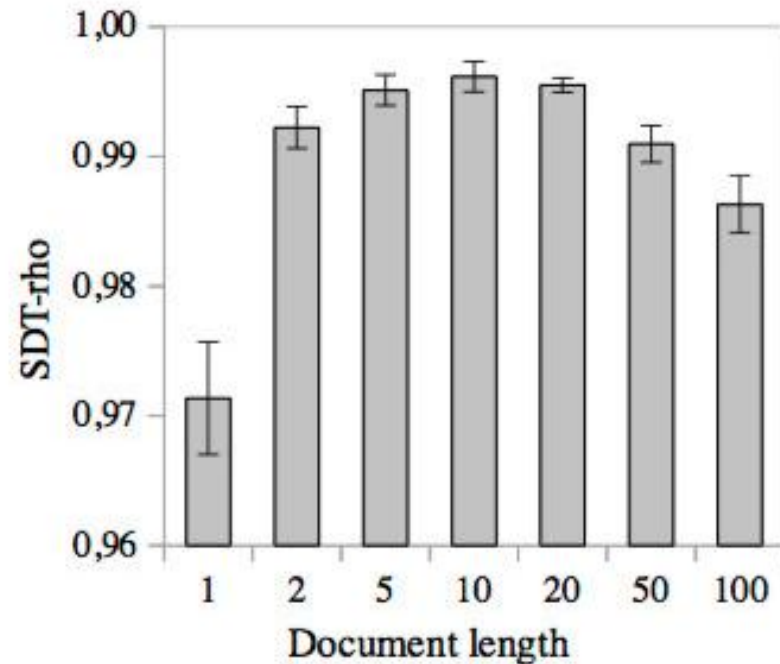
Higher is better

# Results I

Correlation :  
Document length



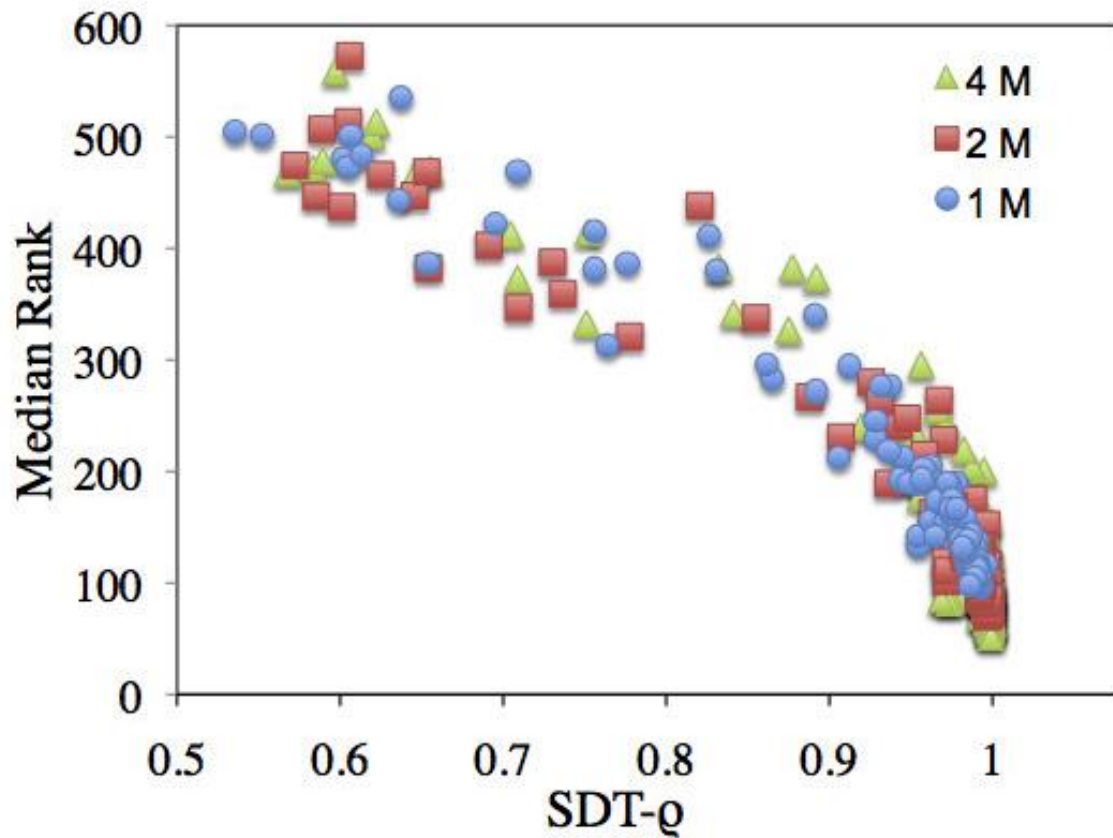
Lower is better



Higher is better

# Results I

Correlation :  
Overall



MIC=0.677  
p<0.0001

# Experiment II

---

TOEFL Synonym test ( Landauer & Dumais, 1997)

Costly

→ beautiful

→ expensive

→ popular

→ complicated

Prolific

→ productive

→ serious

→ capable

→ promising

.....

% of Correct Answers

?

SDT- $\rho$





# Evaluation Method

---

SDT- $\rho$



→ Intrinsic : no external resources or mastery of the language required



→ Predict human judgement

→ Easy to implement

# How to use in practice

- Correlation computed using the overlap with the benchmarks (Norms, TOEFL)
- What set of words to use as Pseudo-Synonyms in computing the SDT- $\rho$ ?
  - Random words?
  - Set size?
  - Word frequency?

## Correlation with Median Rank

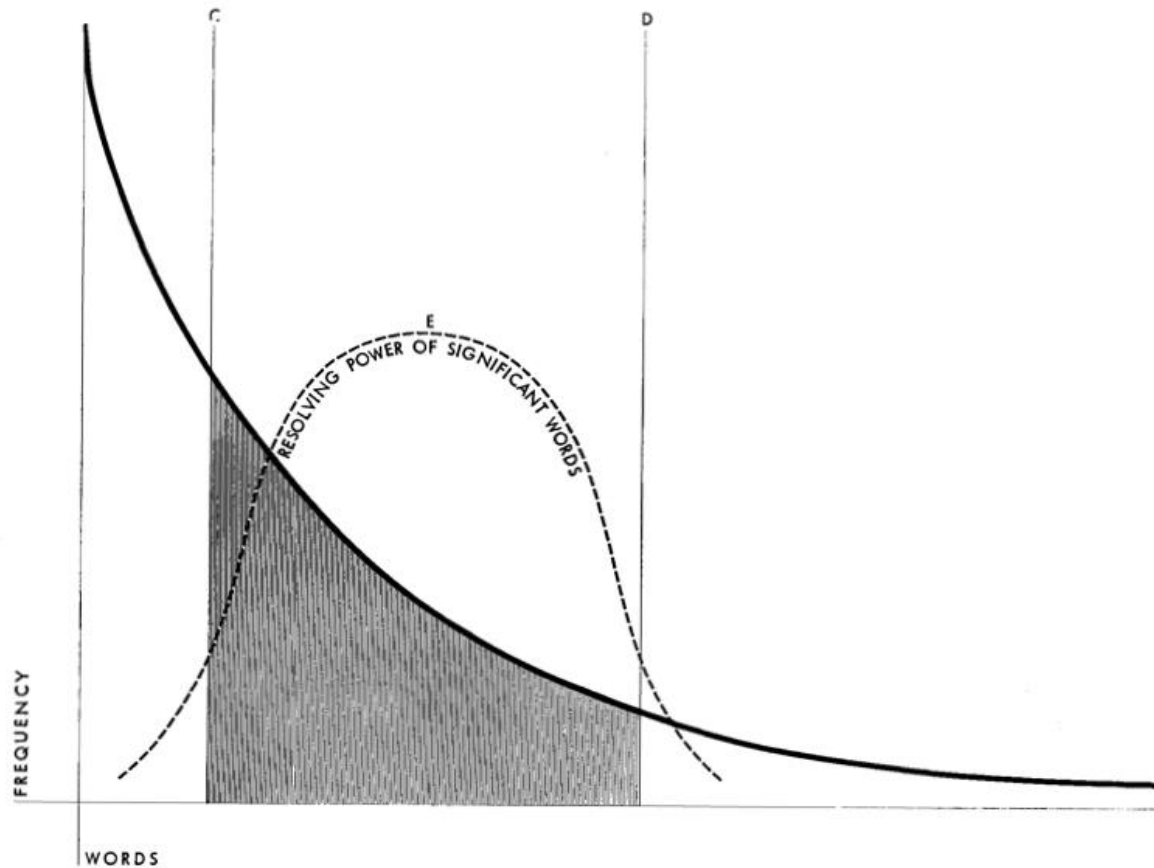
Freq. $x$	$1 < x < 40$			$40 < x < 400$			$x > 400$			All	Overlap
	100	500	1000	100	500	1000	100	500	1000		
Size	100	500	1000	100	500	1000	100	500	1000	~ 4 M	1093
MIC	0.311	0.219	0.549*	0.549*	<b>0.717*</b>	<b>0.717*</b>	0.311	0.205	0.419	0.549*	<b>0.717*</b>

\* :  $p < 0.05$

A small set of random mid-frequency words give the highest correlation with human judgement

# How to use in practice

Mid-frequency words have high discriminating  
(« resolving ») power



Luhn (1958)

# Evaluation Method

---

SDT- $\rho$



→ Intrinsic : no external resources or mastery of the language required



→ Predict human judgement



→ Easy to implement

# Conclusion

---

→ An evaluation method to be used as a proxy when no human generated data is available

→ This method could be used to set the parameters of the semantic models(semantic dimension, doc length, corpus size,..)

→ Gives a global and rough measure of the quality of the semantic structure

→ **Future work** : a word level variant of the method will, a priori, enable us to assess some fine grained properties (distribution properties of abstract/concrete words or linguistic categories,...)

# Questions

---

Thanks for your attention!

Questions?