# Search Powered by Deep Learning

*- From Content Similarity to Semantic Similarity*

Infosys® | Building Tomorrow's Enterprise

# Speakers

John Kuriakose
Principal Architect
Big Data & Analytics - IIP

Ashish Vishwas Kaduskar
Architect
Big Data & Analytics - IIP

Debanjan Mahata
Senior Research Associate
Big Data & Analytics - IIP

# Infosys Information Platform (IIP)



**IIP layers**

- Open Source / Spark Components
- ETL / Integration
- Spark / Storm / Others
- HIVE / HBase/ GraphX / Others
- Hadoop / FS Storage / Infra Management

**Infosys & Partner IP Components**
Tools | Data Extractors | Algorithms | Packaging & Support

**Customization, Integration & Implementation Services**
Data Modeling & Cleansing | Agile App Development
Data Science & Analytics | Security & Governance
Custom Data Extractors

Infosys® | Building Tomorrow's Enterprise

# Agenda

- **What are Word and Doc Embeddings ?**

- **How do they enrich Search Applications ?**

- **How to build them for a search application ?**

- **How to integrate them in a search application ?**

# Searching across Research Articles

recurrent neural networks 🔍

1. **Recognizing recurrent neural networks (rRNN): Bayesian Inference for recurrent neural networks**

   ["Bitzer, Sebastian","Kiebel, Stefan J."] - Fri Jan 20 00:00:00 UTC 2012

   ["recurrent neural networks","bayesian inference","computational neuroscience","machine learning applications","rnn","nonlinear function","brain...

   Recurrent neural networks (RNNs) are widely used in computational neuroscience and machine learning applications. In an RNN, each neuron computes its output as a nonlinear function of its integrated input. While the...

2. **Conversion of Artificial Recurrent Neural Networks to Spiking Neural Networks for Low-power Neuromorphic Hardware**

   ["Diehl, Peter U.","Zarrella, Guido","Cassidy, Andrew","Pedroni, Bruno U.","Neftci, Emre"] - Sat Jan 16 00:00:00 UTC 2016

   ["rnn","artificial recurrent neural networks","low-power neuromorphic hardware","neuromorphic low-power systems","significant momentum","recurrent...

   In recent years the field of neuromorphic low-power systems that consume orders of magnitude less power gained significant momentum. However, their wider use is still hindered by the lack of algorithms that can harness the...

3. **Sequence Modeling using Gated Recurrent Neural Networks**

   Pezeshki, Mohammad - Thu Jan 01 00:00:00 UTC 2015

   ["gated recurrent neural networks","recurrent neural networks","human motion data","next immediate data point","recently proposed gated recurrent units","promisi...

   In this paper, we have used Recurrent Neural Networks to capture and model

**Ingestion**

- Data Preprocessing
- Indexing

**Enrichment**

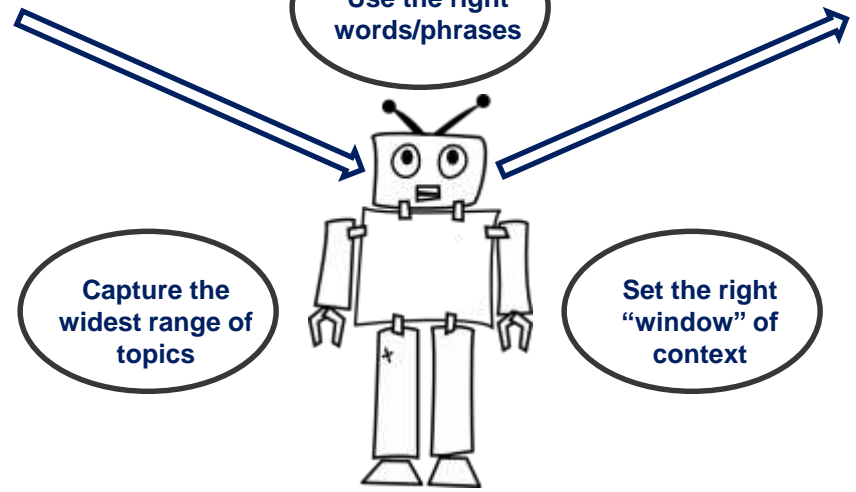- Training different word embedding models

**Content Selection**

- Query Expansion
- Similar Article Recommendation
- Keyword Extraction
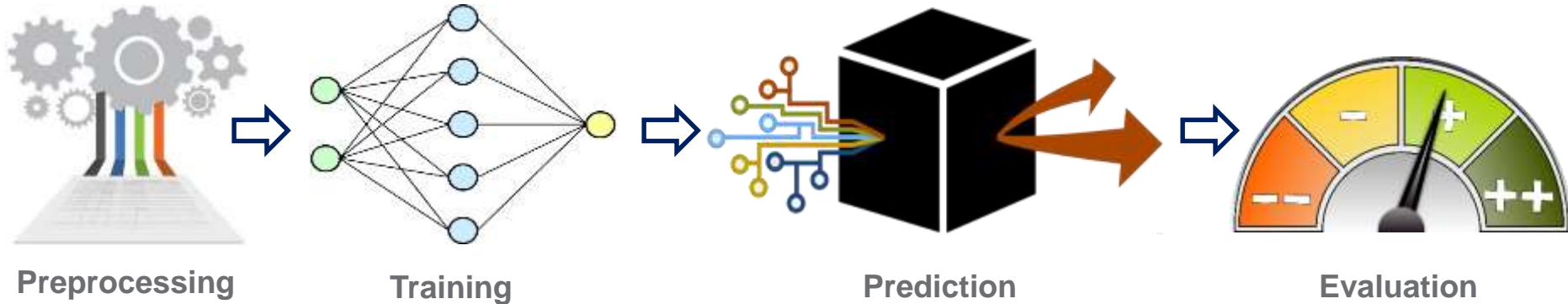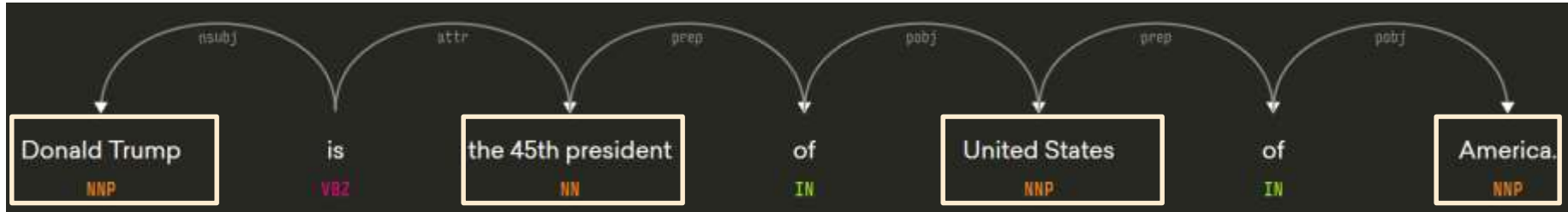
Infosys® | Building Tomorrow's Enterprise

# Humans Vs Machines

**1147001 Research Abstracts**

**Use the right words/phrases**

**Capture the widest range of topics**

**Set the right "window" of context**

# Building the Models



Preprocessing     Training     Prediction     Evaluation

**gensim**
topic modelling for humans

Infosys® | Building Tomorrow's Enterprise

# Preprocessing



**Sentence splitting**

**Phrase Tokenization**

**Removal of strings containing only numeric characters**

**Removal of functional words like 'accordingly', 'although', etc**

**Removal of named entities of types "DATE", "TIME", "PERCENT", "MONEY", "QUANTITY", "ORDINAL", "CARDINAL"**

spaCy

Infosys® | Building Tomorrow's Enterprise

# Word Embeddings

Mikolov et al, (2013a).

# Semantic Similarity between Words



The
Sicilian
gelato
was
extremely
rich.

Sicilian
Italian

extremely
very

ice-cream
gelato

velvety
rich

The
Italian
ice-cream
was
very
velvety.

**word2vec embedding**

Kusner et al (2015)

Infosys® | Building Tomorrow's Enterprise

# Word2Vec



CBOW                    Skipgram

Mikolov et al, (2013a).

# Word2Vec Training



| Source Text | Training Samples |
|---|---|
| **The** quick brown fox jumps over the lazy dog. ⟹ | (the, quick)<br>(the, brown) |
| The **quick** brown fox jumps over the lazy dog. ⟹ | (quick, the)<br>(quick, brown)<br>(quick, fox) |
| The quick **brown** fox jumps over the lazy dog. ⟹ | (brown, the)<br>(brown, quick)<br>(brown, fox)<br>(brown, jumps) |
| The quick brown **fox** jumps over the lazy dog. ⟹ | (fox, quick)<br>(fox, brown)<br>(fox, jumps)<br>(fox, over) |

Source: http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/

# Word2Vec Training

# Word2Vec Training



Hidden Layer Weight Matrix → Word Vector Lookup Table!

300 neurons

300 features

10,000 words

10,000 words

Source: http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/
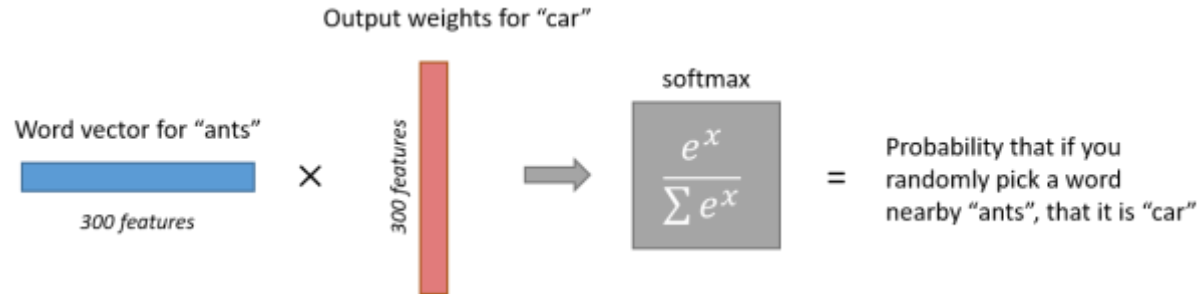
Infosys® | Building Tomorrow's Enterprise

# Word2Vec Training

$$[0 \quad 0 \quad 0 \quad 1 \quad 0] \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = [10 \quad 12 \quad 19]$$

Output weights for "car"

Word vector for "ants"

300 features

$\times$

300 features

softmax

$$\frac{e^x}{\sum e^x}$$

= Probability that if you randomly pick a word nearby "ants", that it is "car"

Negative Sampling

$$log\,\sigma(\nu'_{w0}{}^{T}\nu_{wI}) + \sum_{i=1}^{k} \mathbb{E}_{wi} \sim P_n(w)[log\,\sigma(\nu'_{w0}{}^{T}\nu_{wI})]$$

$$P_n = U(s)^{3/4}/Z$$

Source: http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/
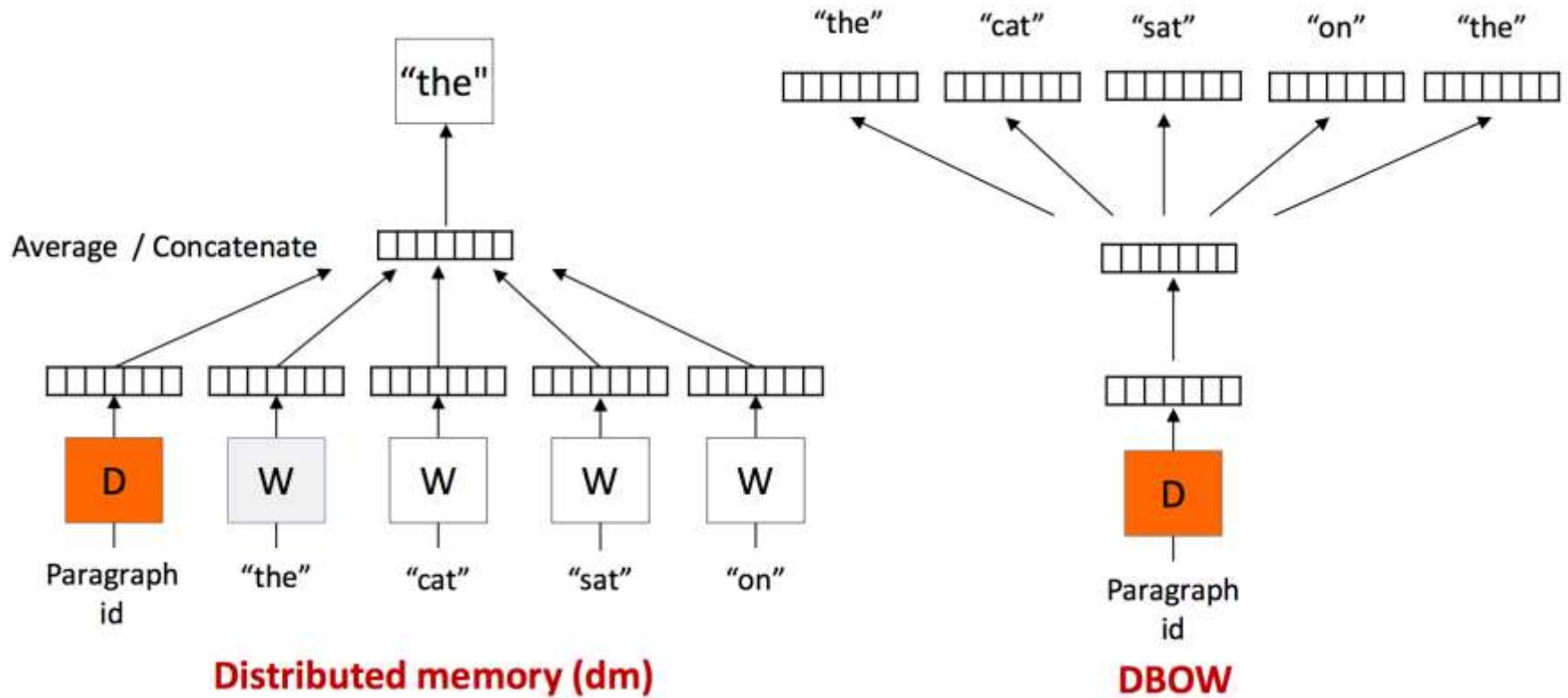
Infosys® | Building Tomorrow's Enterprise

# Fasttext

- Very Similar to Word2vec

- Primary Difference

  - Takes into account the internal structure of words while learning word representations

  - The resultant word vector is a combination of vectors for its constituent character ngrams

- Extremely fast training

- Very good for morphologically rich languages

- Takes into account both "Semantic as well as Syntactic Similarity"

- Very good performances in Syntactic Similarity tasks

# Doc2Vec



Mikolov et al. (2014)

# Training Parameters

| Word2vec | Fasttext | Doc2Vec |
|----------|----------|---------|
| • **Skipgram** | • **Skipgram** | • **Distributed Memory** |
| • **Negative Sampling** | • **Negative Sampling** | • **Dimensions =** 1000 |
| • **Dimensions =** 1000 | • **Dimensions =** 1000 | • **Window Size =** 10 |
| • **Context Window Size =** 5 | • **Max Length of Char Ngrams =** 6 | • **Epochs =** 10 |
| • **Learning Rate =** 0.025 | • **Min Length of Char Ngrams =** 3 | • **Initial Learning Rate =** 0.025 |
| • **Trained on unigrams and phrases** | • **Learning Rate =** 0.05 | • **Trained on unigrams and phrases** |
| | • **Trained on unigrams and phrases** | |

Infosys® | Building Tomorrow's Enterprise

# Query Expansion

natural language processing

natural language generation

information extraction

text categorization          nlp          🔍          machine translation

text mining

knowledge discovery

word sense disambiguation

- **Getting rid of thesaurus based or dictionary based query expansion.**
- **How many phrases to use for expansion ?**
- **Which model ?**
- **How to integrate the trained models with a search application  for query expansion ?**

# Recommending Similar Research Articles



**Content Similarity**

**Vs**

**Semantic Similarity**

# Ranked Keyphrase Extraction

| Supervised | Unsupervised |
|---|---|
| • **Known to give better results**<br><br>• **Drawbacks**<br><br>  – **Domain Specific**<br><br>  – **Training and Tuning of the models for generalization**<br><br>  – **Intelligent Feature Engineering**<br><br>• **Examples: KEA, MAUI** | • **Known to give worse results than supervised in domain specific tasks**<br><br>• **Domain independent**<br><br>• **No need of feature engineering**<br><br>• **Algorithm determines relationships between candidates for identifying the keyphrases**<br><br>• **Examples: TextRank, RAKE** |

# Ranked Keyphrase Extraction

# Integration

## Ingestion

**Indexing**

**Preprocessing**

## Content Selection

**Query expansion**

**Document Recommendation**

**Ranked Keyword Extraction**

**API**

## Enrichment

**Word Embedding Model**

**Doc Embedding Model**

**Keyword classifier**

# References

1. Sebastian Rudder, "**On Word Embeddings**"; http://sebastianruder.com/tag/word-embeddings/index.html

2. Christopher Olah, "**Colah's Blog**", http://colah.github.io/

3. Chris McCormick, "**Word2Vec Tutorial – The Skip-Gram Model**", http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/

4. Le, Quoc V., and Tomas Mikolov. "**Distributed Representations of Sentences and Documents.**" ICML. Vol. 14. 2014.

5. Mikolov, Tomas, et al. "**Distributed representations of words and phrases and their compositionality.**" Advances in neural information processing systems. 2013.

6. Bojanowski, Piotr, et al. "**Enriching word vectors with subword information.**" arXiv preprint arXiv:1607.04606 (2016).

7. Rong, Xin. "**word2vec parameter learning explained**." arXiv preprint arXiv:1411.2738 (2014).

Infosys® | **Building Tomorrow's** Enterprise