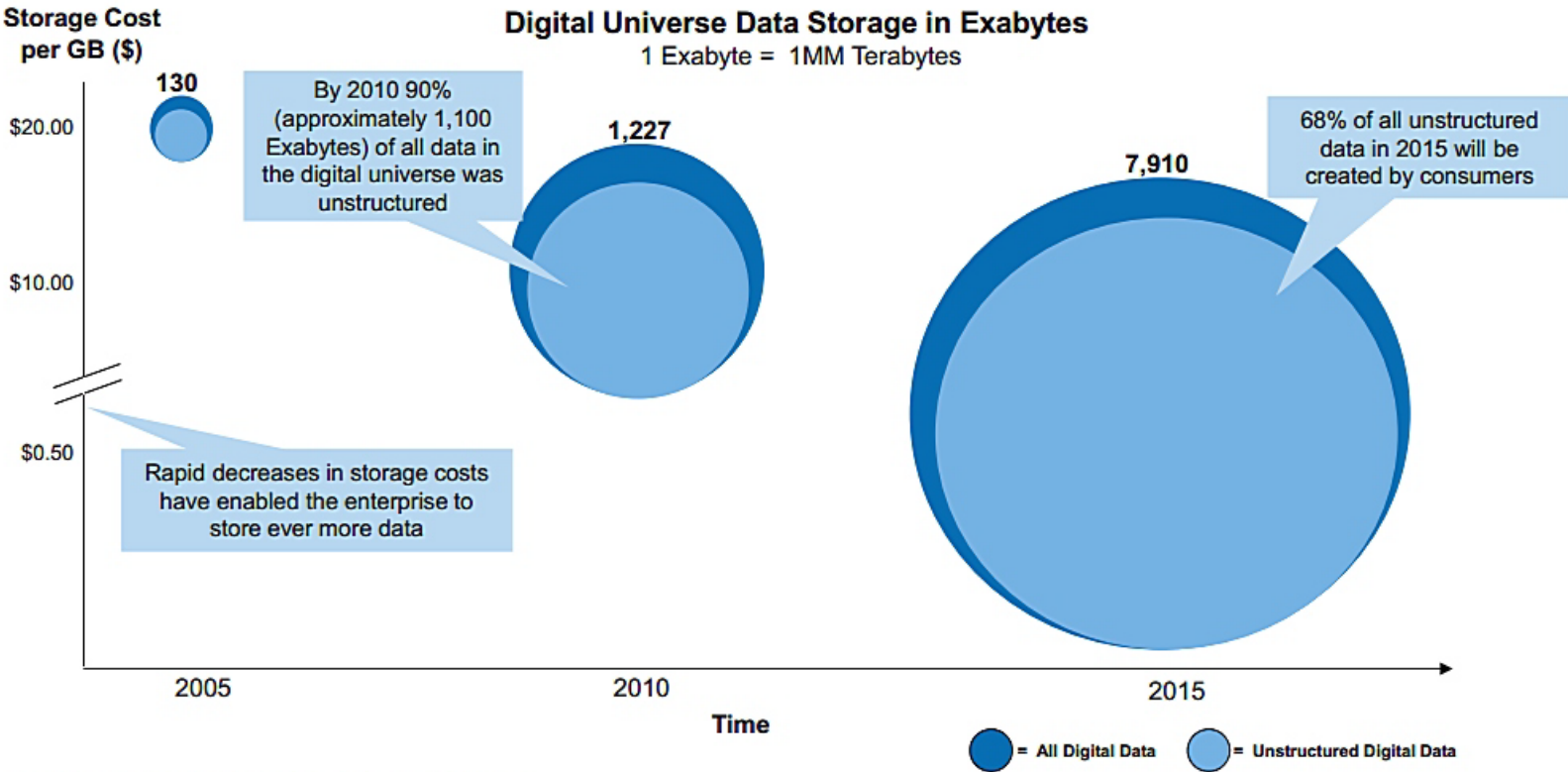# Distributed Representation of Documents with Explicit Explanatory Features

Han Kyul Kim, Hyunjoong Kim, Sungzoon Cho

Dep. of Industrial Engineering

Seoul National University

# Ascendance of "Unstructured Data"

- Due to the proliferation of smart phones and decreasing data storage cost, the amount of data being generated has drastically increased over the years

- Unstructured data is a realm of unexplored potential that provides an insight into consumers' psyche, as most of them are being generated directly from the consumers,
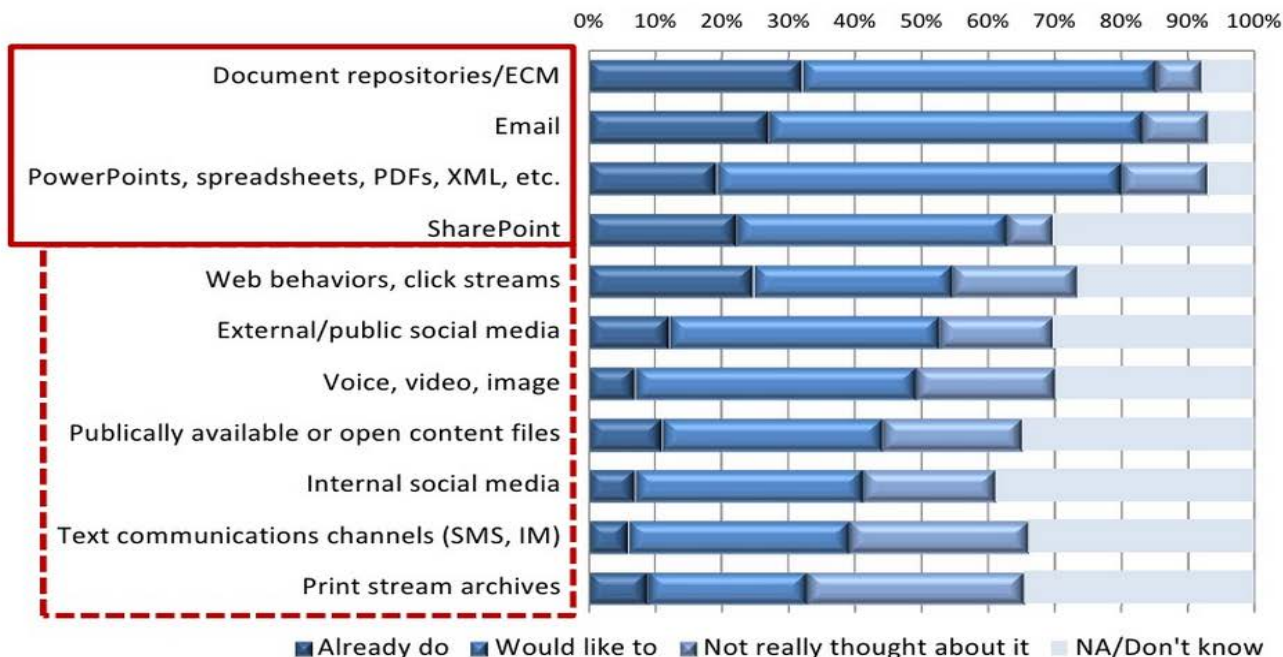


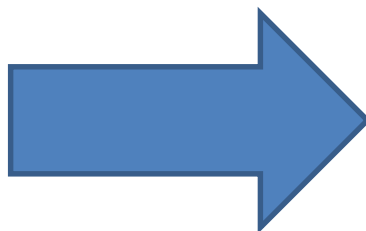Source: IDC's Digital Universe Study, sponsored by EMC, June 2011

- Text data, as one of the most common form of unstructured data, has generated great interest amongst companies

- Extracting insight through text mining essentially begins with analyzing documents



Are there large unstructured or semi-structured data repositories (ie, text, rich media, etc.) in your business that you would like to analyze, monitor or query - as opposed to search/retrieve?

©AIIM 2012          N=300

- Words
- Paragraphs
- Themes
- Main arguments

......

$Doc\ 1: [1, 4, 3, 0, … ]$
$Doc\ 2: [5, 2, 8\ 0, … ]$
.........

- A document is a complex system composed with words, paragraphs and many explicit and implicit features

- In order for computers to extract insights from such complex system, documents need to be represented in some kind of numerical format

# Bag of Words Approach

- <u>Bag of Word Hypothesis:</u> frequencies of words in a document tend to indicate the relevance of the documents

- Basically, BOW uses token frequencies as features for representing a document

**[Document 1]:**

Arsenal's annual injury problem is underway. Their thin squad will be put to the test by a Swansea City team looking to build on a vital win at Aston Villa.

Arsenal have a whole host of injury problems to contend with. The Gunners currently sit top of the Premier League's infamous injury table. Eight senior players will be unable to take part at the Liberty Stadium

|  | **Arsenal's** | **Annual** | **Injury** | **Problem** | **is** | **...** |
|---|---|---|---|---|---|---|
| **Document 1** | 1 | 1 | 3 | 1 | 1 | ... |

# Bag of Words Approach

- Pros: It's intuitive!
    - A feature of document vectors uniquely correspond to a specific token
    - Users can directly understand the components of the documents
    - It can provide clear explanation for why and how certain documents are different from each other

- Cons: Dimension of a vector increases quickly
    - A number of words in a document is NOT small
    - Words with similar meanings (not necessarily synonyms) and identical words with different grammatical structures are considered as individual tokens

**[Document 1]:**

**Arsenal's** annual injury **problem** is underway. **Their** thin squad will be put to the test by a Swansea City team looking to build on a vital win at Aston Villa.
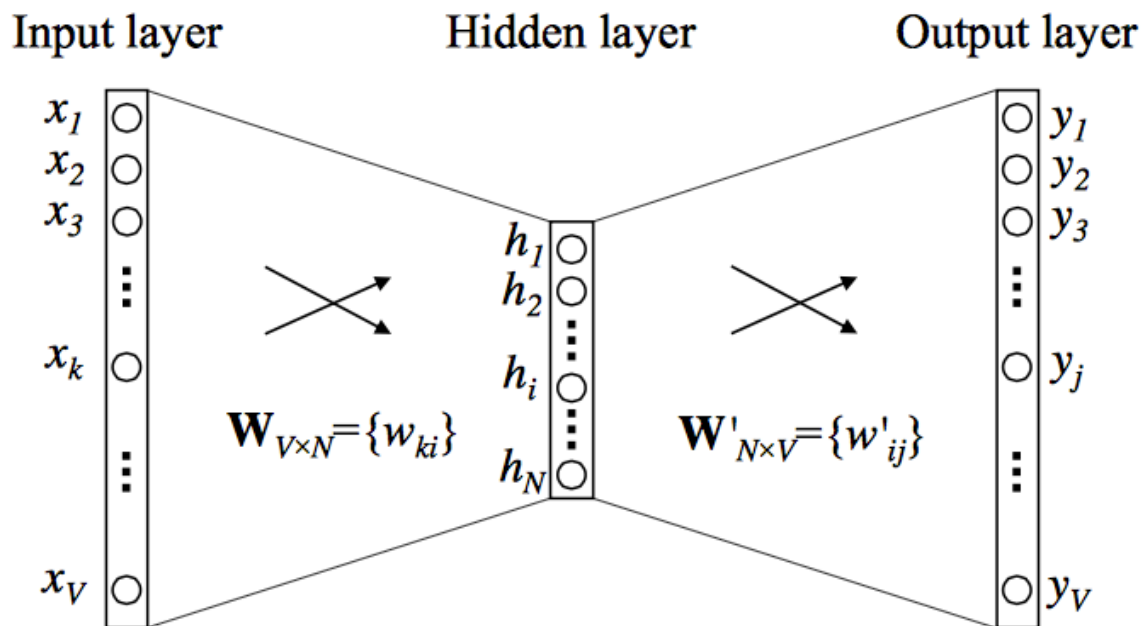
**Arsenal** have a whole host of injury **problems** to contend with. **The Gunners** currently sit top of the Premier League's infamous injury table. Eight senior players will be unable to take part at the Liberty Stadium

## Distributed Representation: Word2Vec & Doc2Vec

**<u>Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." Advances in Neural Information Processing Systems. 2013.</u>**

- <u>Distributed Hypothesis</u>: words that occur in similar contexts (neighboring words) tend to have similar meanings

- Based on this assumption, simple neural network is used to embed words into continuous vector space

- Train the weights of the network so that an input word to the network can predict its neighboring words within certain window size

# Simple Word2Vec Architecture

- Amongst vocabulary of size V, let's say we want to predict one target word(output) when we are given one context word (input) → bigram structure

- Input vector is an one-hot encoded vector (only one node with for designated context word will be 1)

- Weight matrix $W_{VXN}$ and $W'_{NXV}$ embed context and word into continuous vector space, respectively
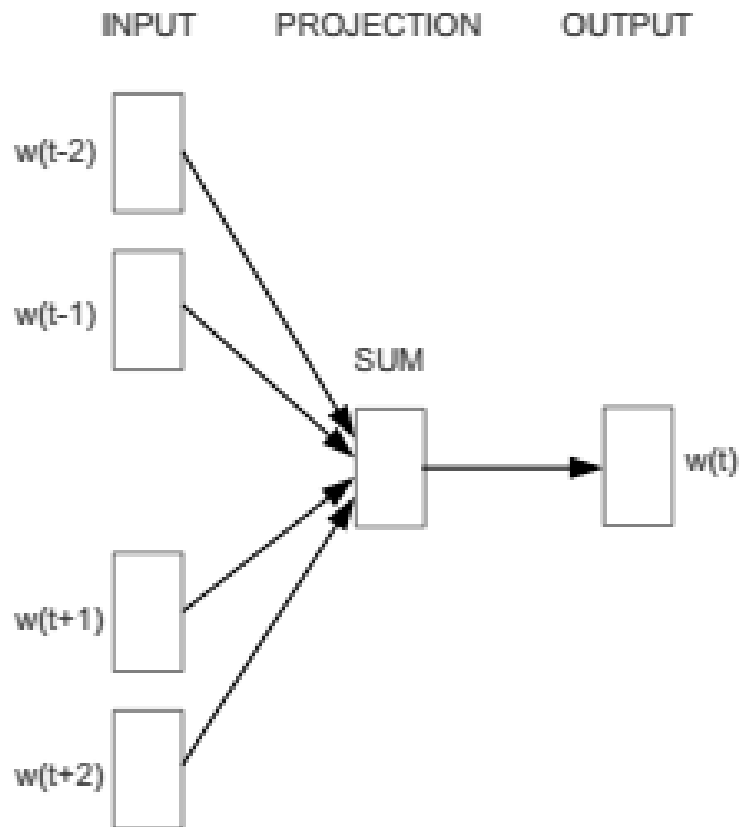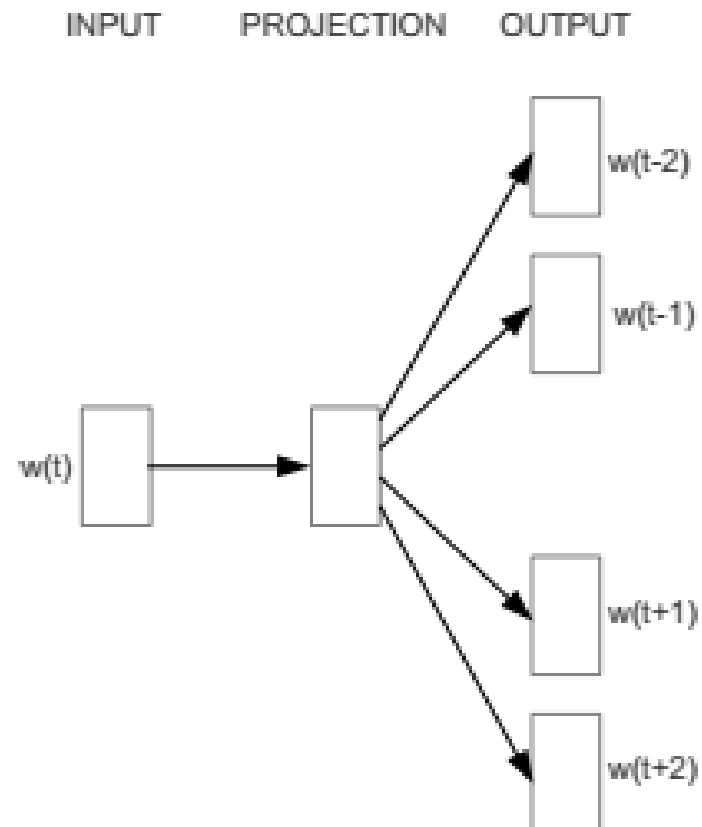
# Simple Word2Vec Architecture

- Non-linearity (soft-max) is incorporated only for computing the output

- $v_{w_I}$: vector representation of the input context word $w_I$ ($x^T W$, k-th row of W)

- $v'_{w'_j}$: vector representation of the output word (j-th column of W')

- Train the weights via gradient descent to maximize conditional (log) probability of observing the actual output word $w_O$ given the input context word $w_I$

- As a result, word vectors with similar context (thus, meaning) will be located close to each other in the embedding space

$$p(w_j|w_I) = \frac{\exp\left(\mathbf{v'_{w_O}}^T \mathbf{v}_{w_I}\right)}{\sum_{j'=1}^{V} \exp\left(\mathbf{v'_{w'_j}}^T \mathbf{v}_{w_I}\right)}$$

**Example**

China
Korea
Japan
America

Chinese
American
Korean
Japanese

Beijing
Seoul
Tokyo
D. C.

Mommy
Daddy
Mom
Dad
Uncle
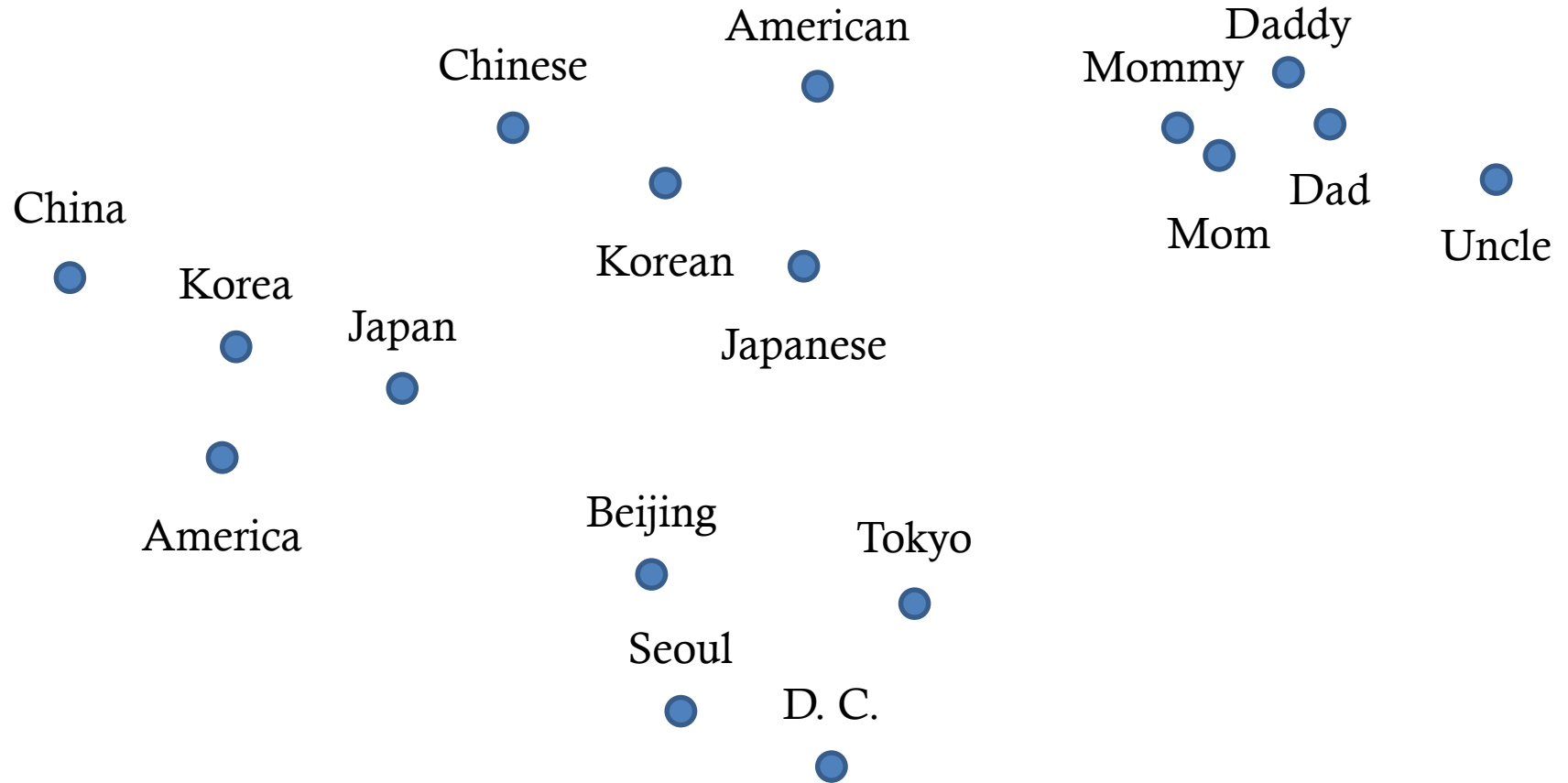
- Word vectors with similar context (thus, belonging to a similar concept) will be located close to each other in the resulting embedding space

**Xing, Chao, et al. "Document classification with distributions of word vectors." Asia-Pacific Signal and Information Processing Association, 2014 Annual Summit and Conference (APSIPA). IEEE, 2014.**

- **Simple average pooling approach:**

$$v_i = \frac{1}{J_i} \sum_{j=1}^{J_i} c_{i,j}$$

- Derives a document vector as the centroid of word vectors within the document

- Bias towards words without significant contribution to representing the semantics of the documents

<Word2Vec>                    <Paragraph2Vec>

**Le, Q. V., & Mikolov, T. (2014). Distributed representations of sentences and documents**

- Extension of Word2Vec: a document is considered as an extra word

- Document (paragraph) id represents one-hot encoded vector of <u>documents</u>

- As a result, documents are also embedded into continuous vector space

Figure 4: Results of experiments on the hand-built Wikipedia triplet dataset.

Table 3: Performances of different methods on hand-built triplets of Wikipedia articles on the best performing dimensionality.

| Model | Embedding dimensions/topics | Accuracy |
|---|---|---|
| Paragraph vectors | 10000 | 93.0% |
| LDA | 5000 | 82% |
| Averaged word embeddings | 3000 | 84.9% |
| Bag of words | | 86.0% |

**Dai, Andrew M., et al. "Document Embedding with Paragraph Vectors." NIPS Deep Learning Workshop. 2014.**

13

# Distributed Representation Approach

- <u>Pros:</u> Good Performance!
  - Document vector successfully captures useful yet unknown features for representing documents (validated through document clustering and classification)
  - Dimension of a document vector is restricted to certain size

- <u>Cons:</u> Not intuitive!
  - As with other neural network based models, unknown feature are used to represent a document
  - Each value in a document vector doesn't provide any explicit explanation about the document

[-0.08759557455778122, -0.04312118515372276, -0.08494572341442108, 0.024585919454693794, -0.05785191431641579, -0.02659076638519764, 0.04704275727272034, -0.03940117731690407, 0.005195754114538431, -0.018994472920894623, -0.03089658915965515, -0.02599106915295124, -0.029802896082401276, -0.009517285041511059, -0.03624524921178818, 0.0029738633893430233, -0.04270448908209801, -0.0890769511461258, -0.04064304754137993, 0.017775749787688255, 0.0910411849617958, 0.05333533510565758, -0.07692492008209229, 0.08628936856985092, -0.042326122522354126, -0.007681592833250761, 0.0414172001183033, -0.03035894976156044, 0.05717118829488754, 0.0396726056933403, -0.09482061862945557, 0.0538295470178127, -0.016189705580472946, 0.0013550696894526482, 0.00425155786797404, -0.10439810156822205, -0.01734139770269394, 0.0873356834053932, -0.020141418437245, 0.06905293464660645, -0.05219369381661835, -0.008379205130040646, 0.050789929926395416, -0.0521097406744957, 0.02524719014763832, -0.09064795076847076, -0.01605154387652874, -0.08548879623413086, 0.09579522907733917, 0.07222563773393631, -0.01747663877904415, -0.07119490951299667, 0.04312814772129059, 0.006512957159429789, 0.04662078991532326, 0.053695641458034515, 0.0017072528135031462, 0.0468018501996994, 0.03436211124062538, -0.0556086413562978, 0.04996718838810921, 0.09512156248092651, 0.0048730792477269, -0.0013316450640559196, 0.063605397939682, 0.02428655326366424, 0.0097620347514748, -0.04571126028895378, 0.03228874504561926, -0.06711595505475998, 0.03179262951016426, 0.00038015464087948203, 0.013781636022031307, 0.049716588109731674, -0.04940832788811455, 0.0465905033051967, 0.08707352727651596, -0.10198234766721725, 0.0012964112684130669, 0.019826047122478485, 0.02079876139760017, 0.03980609402060509, -0.016105623915791510, 0.09880314767360687, 0.035302355885505676, -0.03354038670659065, -0.060332611203193665, 0.009992017410695553, -0.07922962307929993, -0.04672875627875328, 0.02924907766282558, 0.0073861079290509, 0.01445483043789863, 0.0367921367287635, 0.0225952658802270, -0.07544630010890, 0.0375224202871322, -0.042987767606973, 0.0429357662796974, -0.04286639392375946, 0.0549306534230709, -0.010759958997368813, -0.0264284871518611, 0.0593683943152427, 0.0107355564837458, -0.02210656180977821, 0.00509836338460445, 0.02817825973033905, 0.06781460344791412, 0.01152470614761141, 0.04529837146401405, -0.1041417792439460, -0.06333499401807785, -0.025369135662913322, 0.01380544807761907, 0.088760338723565952, -0.03237186372280121, 0.08923118561506271, 0.035788971930742264, -0.073584914207458, -0.114317052066326, 0.004165078978985548, 0.03989437595009804, 0.0116985784843564, 0.000912736915051937, -0.0071949455887079, -0.06279811263084412, -0.01203678268939256, 0.040613092482089996, -0.07241667807102203, 0.068035975098600992, 0.0318059809505939, 0.020111583173274994, 0.05419156327843666, -0.06192755699157715, 0.0306070601427551, -0.015478908084332943, -0.05181756243109703, -0.02790454049804687, 0.0161539483815431, 0.00118022342212498, -0.0524372085928916, -0.05809084326028824, 0.0256935339421033, -0.0445638187229633, -0.0523514263331890, 0.0028492037891674, 0.0251784324645996, 0.0742745026946067, -0.0041514914482831955, 0.067084312438964, 0.01024503447115421, -0.037210747599601746, 0.0519414506852626, 0.011150983126322552309, -0.0136907296253209, -0.00447499845176935, 0.0840763598940962, 0.018420604989326, 0.0693516209721565, -0.07085989415645, -0.04950730875134468, 0.1205481216311454, -0.006143994629383087, 0.01835456631302834, 0.10588039457798004, 0.04850513488054275, -0.014145134948194027, -0.02343590557575226, -0.00773952063173055, -0.033509612083435, 0.0214711278676986, 0.07666371762752533, 0.06256308406591415, 0.01227859035134315, -0.0174938272684812, -0.0258525852113962, 0.0285028610378503, -0.04092925786972046, 0.138949438929557, -0.000527621130459010, 0.019603941589593887, 0.04787249863147735, -0.024796431884169, -0.03278869017958641, -0.03298942744731903, -0.07497345651565, 0.01786929555237293, 0.0143599743023514, 0.0326016172766685, -0.058985892683267, 0.04998610913753509, -0.01653622649610042, -0.0198389366269111, -0.03034561872482299, -0.045678805559873, -0.046361502259969, -0.07739298790693283, 0.0545377209782600, -0.01512290351092815, -0.05582341924309730, 0.0209887269884347, 0.0380556583404541, 0.041010051965713, 0.065005554829835892, -0.01441126875579357, 0.063243038952306, 0.1223570033907890, -0.00313800743790, 0.041506013093376, -0.05469806492328644, 0.0281550716608762, -0.025410173460841, 0.0236208252608776, -0.05574937537312507, 0.03488092869520187, 0.0292510241270065, -0.04492440819740295, -0.05796622857451439, 0.160918012261390, -0.01372791174799203, -0.002823803108185529, 0.00505891163359880, 0.00933037512004375, -0.04079080000519752, 0.0126318996772170, -0.03498981520533562, -0.0907597243785858, -0.0406805910170078, 0.056468620896339417, 0.0793653130531311, -0.10637512058019638, 0.0249171499162912, 0.0766419768333435, 0.

# 1. Train Word2Vec from the collection of documents

## 2. Cluster word2vec generated vectors to create clusters of concepts

3. **Represent the documents by counting the number of times that their words belong to these different concept clusters (Similar to BOW approach!)**

**Concept Cluster 1** = {Arsenal, Arsenal's, Aston Villa, Swansea City, Gunners...}
**Concept Cluster 2** = {Squad, Players...}

**[Document 1]:**
**Arsenal's** annual injury **problem** is underway. **Their** thin squad will be put to the test by a Swansea City team looking to build on a vital win at Aston Villa.

**[Document 2]:**
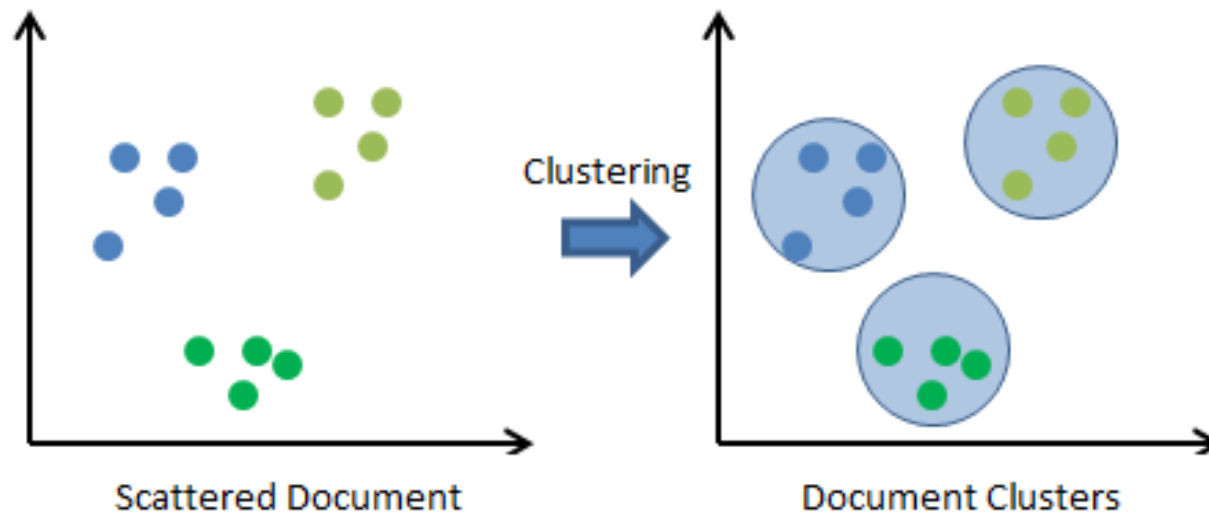**Arsenal** have a whole host of injury **problems** to contend with. **The Gunners** currently sit top of the Premier League's infamous injury table. Eight senior players will be unable to take part at the Liberty Stadium

| Features | Concept Cluster 1 | Concept Cluster 2 | ... |
|----------|:-----------------:|:-----------------:|:---:|
| Document 1 | 3 | 1 | ... |
| Document 2 | 2 | 1 | ... |

4. **Evaluate the effectiveness of the document representation through document clustering task using the document vectors created from the previous step**

| Features | Concept Cluster 1 | Concept Cluster 2 | ... |
|----------|-------------------|-------------------|-----|
| Document 1 | 3 | 1 | ... |
| Document 2 | 2 | 1 | ... |
| ... | .... | .... | ... |
| Document n | 0 | 8 | ... |



Scattered Document → Clustering → Document Clusters

# Dataset: \<Reuters\>

**Total Number of Documents: 203,923 (2006. 09. 01 ~ 2015. 06. 06)**

- Divided into 8 different categories
- Total number of sentences:  3,076,016
- Total number of tokens: 89,146,031
- Total number of unique tokens: 65,159
- Dimension of Doc2Vec and Word2Vec set as 500

| Categories | Number of Documents |
|------------|---------------------|
| Entertainment | 25,500 |
| Sports | 25,500 |
| Technology | 25,500 |
| Market | 25,423 |
| Politics | 25,500 |
| Business | 25,500 |
| World | 25,500 |
| Health | 25,500 |

# Spherical K Means

- To deal with these issues of high dimensional clustering, spherical k means algorithm has been used

- Essentially same as k means algorithm but with cosine similarity as a measure of proximity instead of Euclidean distance

- Chosen best result in terms of inertia

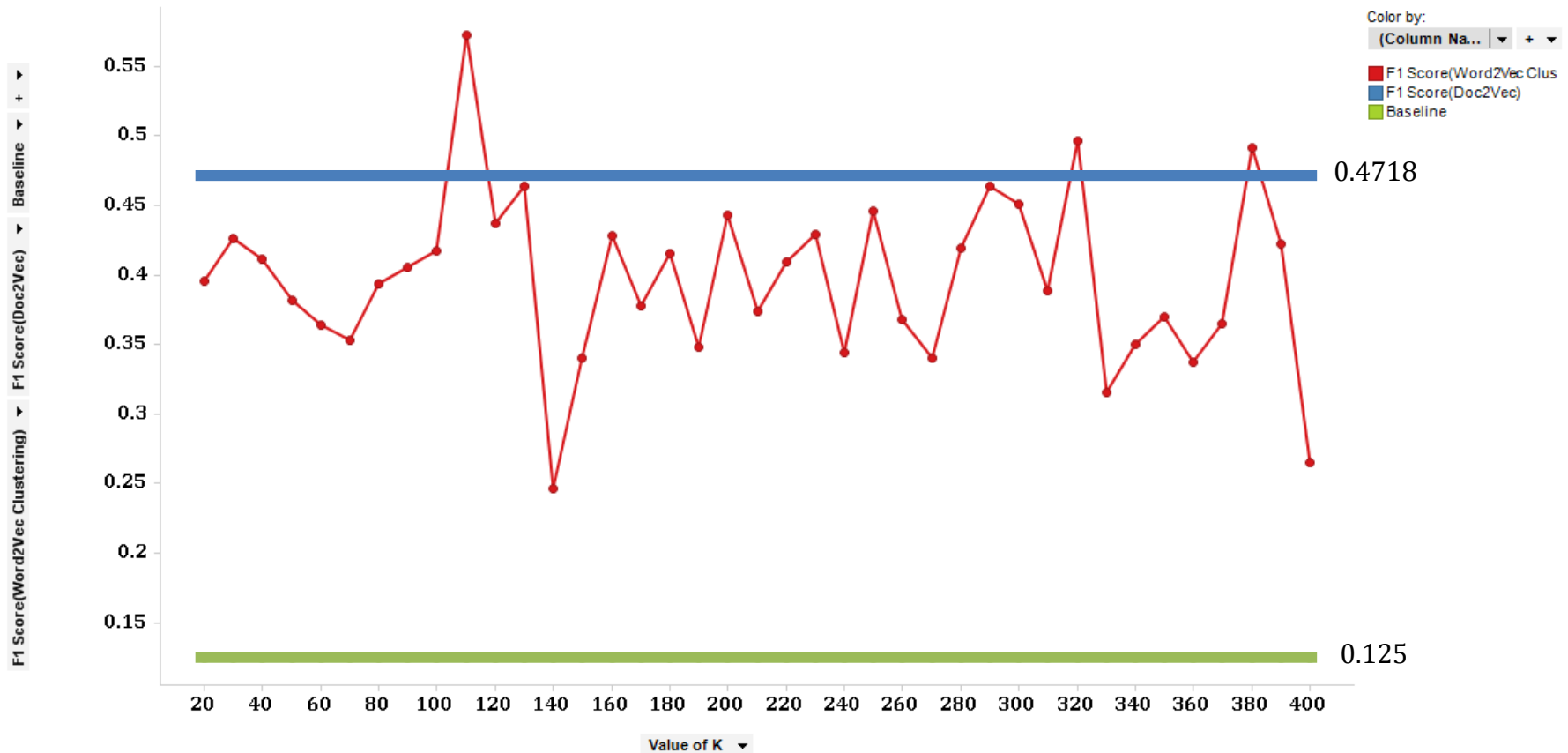**Algorithm:** spherical k-means (SPKM)
**Input:** A set of $N$ *unit-length* data vectors $\mathcal{X} = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$ in $\mathbb{R}^d$ and the number of clusters $K$.
**Output:** A partition of the data vectors given by the cluster identity vector $\mathcal{Y} = \{y_1, ... y_N\}$, $y_n \in \{1, ..., K\}$ .
**Steps:**
1. Initialization: initialize the *unit-length* cluster centroid vectors $\{\mu_1, ..., \mu_K\}$ ;
2. Data assignment: for each data vector $\mathbf{x}_n$, set $y_n = \arg\max_k \mathbf{x}_n^T \mu_k$ ;
3. Centroid estimation: for cluster $k$, let $\mathcal{X}_k = \{\mathbf{x}_n | y_n = k\}$, the centroid is estimated as $\mu_k = \sum_{\mathbf{x} \in \mathcal{X}_k} \mathbf{x} / \| \sum_{\mathbf{x} \in \mathcal{X}_k} \mathbf{x} \|$ ;
4. Stop if $\mathcal{Y}$ does not change, otherwise go back to Step 2a.

Fig. 1.  Spherical k-means algorithm.

- Depending on the value of K used for clustering word vectors, the performance of document vectors vary greatly

- If appropriate value of K is chosen (in this case, K = 110), proposed method can outperform Doc2Vec in document clustering, indicating better document representation

# Word Clusters

- Pharmaceutical Terminologies

| | |
|---|---|
| two-drug | 3 |
| Genzymes | 3 |
| Suppressant | 3 |
| Vectibix | 3 |
| Adderall | 3 |
| Intravenously | 3 |
| Crestor | 3 |
| non-prescription | 3 |
| Antivirals | 3 |
| Sandoz | 3 |
| Discontinuing | 3 |
| anti-retroviral | 3 |
| anti-viral | 3 |

- Military Terminologies

| | |
|---|---|
| dni | 47 |
| Expeditionary | 47 |
| Panettas | 47 |
| maj. | 47 |
| Marshals | 47 |
| Operationally | 47 |
| Chiefs | 47 |
| pentagon. | 47 |
| Commandant | 47 |
| Stationing | 47 |
| Troop | 47 |
| Defences | 47 |
| Lieutenant | 47 |

# Word Clusters

- Garment Related Terminologies

| | |
|---|---|
| handbags | 60 |
| lagerfeld | 60 |
| lipstick | 60 |
| disheveled | 60 |
| flowing | 60 |
| bra | 60 |
| figurines | 60 |
| embellished | 60 |
| hoodie | 60 |
| garlands | 60 |
| masculine | 60 |
| tuxedos | 60 |
| garish | 60 |

- Baseball Terminologies

| | |
|---|---|
| dickey | 72 |
| beltre | 72 |
| peavy | 72 |
| runs | 72 |
| ervin | 72 |
| pinch-hitter | 72 |
| cincinnatis | 72 |
| batting | 72 |
| jacoby | 72 |
| cubs | 72 |
| sabathia | 72 |
| prado | 72 |
| hits | 72 |

# Word Clusters

| Features | X[0] | ... | X[47] | ... | X[72] | ... |
|----------|------|-----|-------|-----|-------|-----|
| Document 1 | 2 | ... | **23** | ... | **0** | ... |
| Document 2 | 2 | ... | **0** | ... | **33** | ... |

Due to 47$^{nd}$ feature, document 1 is something about military.

Due to 72$^{nd}$ feature, document 2 is something about baseball.

# Conclusion

- The proposed method incorporates both the effective performance of the distributed representation and intuitive explanatory power of BOW representation

- For certain value of K used for clustering word2vec vectors, the proposed method, at least in terms of document clustering, can represent the documents at a similar performance level as Doc2Vec

- As the performance of Doc2Vec can be sensitive to the number of dimensions (the number of hidden nodes), the proposed method needs to be compared with various Doc2Vec models resulting from different number of dimensions

- Performance of the proposed method should also be checked for document classification task in order to substantiate the effective document representation of the proposed method

# Reference

Aggarwal, Charu C., Alexander Hinneburg, and Daniel A. Keim. On the surprising behavior of distance metrics in high dimensional space. Springer Berlin Heidelberg, 2001.

Dai, Andrew M., et al. "Document Embedding with Paragraph Vectors." NIPS Deep Learning Workshop. 2014.

Kriegel, Hans-Peter, Peer Kröger, and Arthur Zimek. "Clustering high-dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering." ACM Transactions on Knowledge Discovery from Data (TKDD) 3.1 (2009): 1.

Le, Quoc V., and Tomas Mikolov. "Distributed representations of sentences and documents." arXiv preprint arXiv:1405.4053 (2014).

Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." Advances in Neural Information Processing Systems. 2013.

R. Liu, D. Wang, and C. Xing, "Document classification based on word vectors." ISCSLP, 2014

Rong, Xin. "word2vec Parameter Learning Explained." arXiv preprint arXiv:1411.2738 (2014).
5. IEEE, 2005.

# Reference

Turney, Peter D., and Patrick Pantel. "From frequency to meaning: Vector space models of semantics." Journal of artificial intelligence research 37.1 (2010): 141-188.

Xing, Chao, et al. "Document classification with distributions of word vectors." Asia-Pacific Signal and Information Processing Association, 2014 Annual Summit and Conference (APSIPA). IEEE, 2014.

Zhong, Shi. "Efficient online spherical k-means clustering." Neural Networks, 2005. IJCNN'05. Proceedings. 2005 IEEE International Joint Conference on. Vol. 5. IEEE, 2005.