# Distributed Representation of Documents with Explicit Explanatory Features: Pilot Test
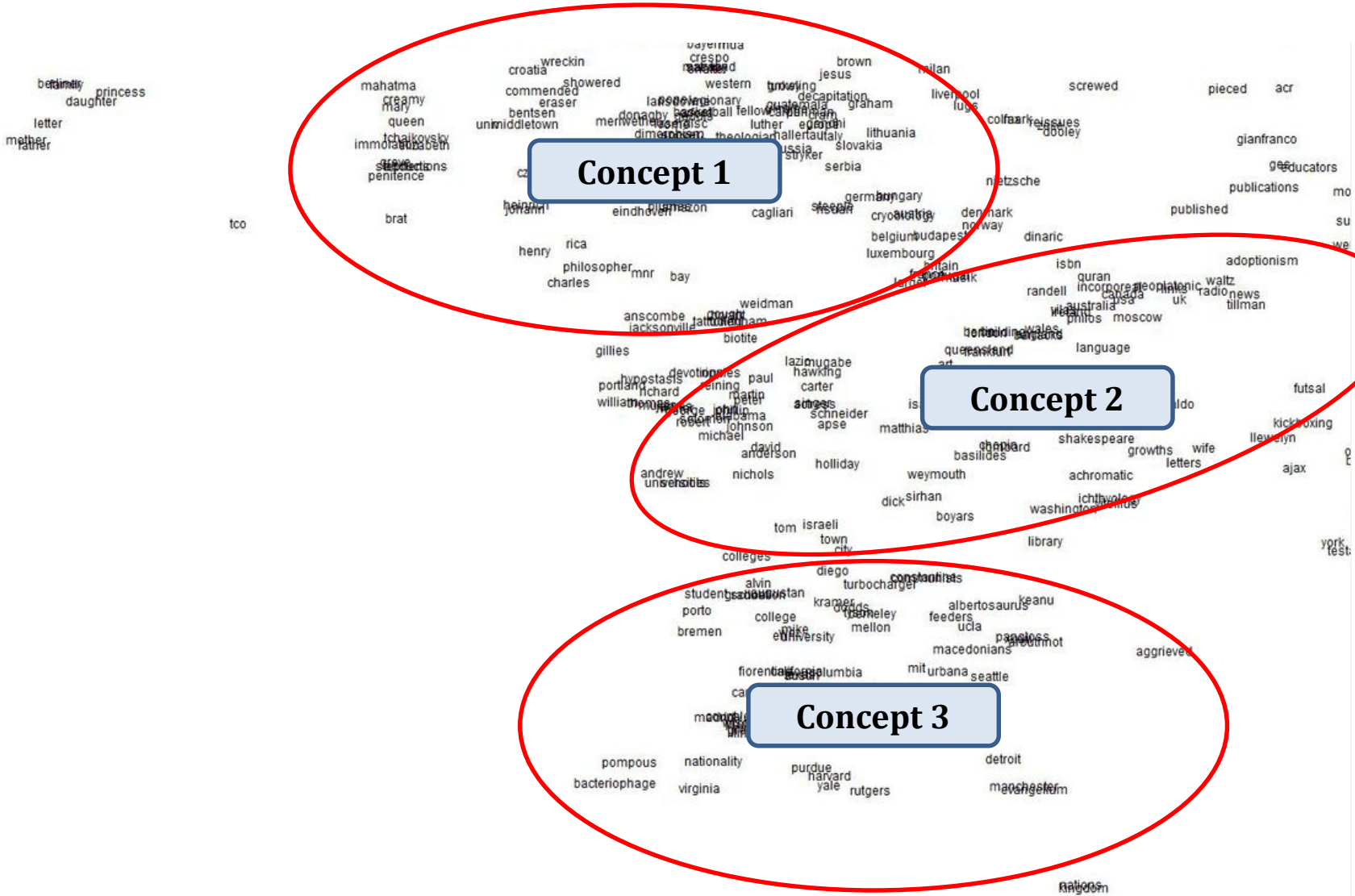
October 7th, 2015
SNU Data Mining Center
Han Kyul Kim

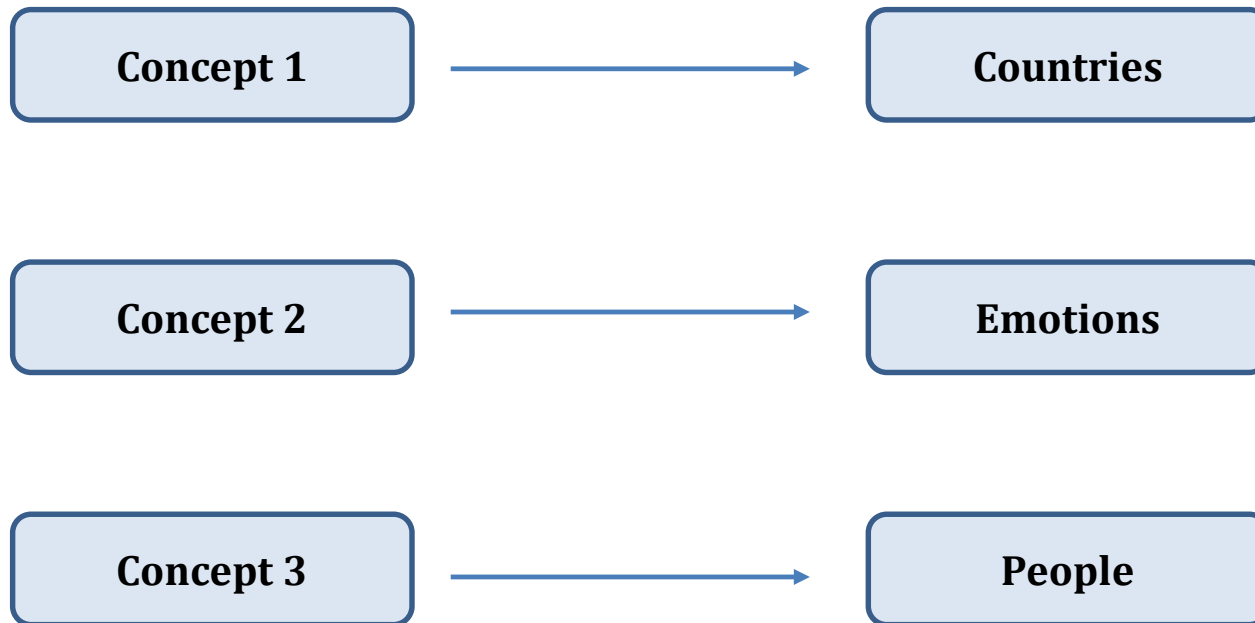1. **Train Word2Vec with the collection of documents**

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

3. Label the concepts using the words associated with each cluster

| Concept 1 | ⟶ | Countries |
| Concept 2 | ⟶ | Emotions |
| Concept 3 | ⟶ | People |

## 4. Represent the documents using the counts of these concepts

Doc 1 = [Countries, Emotions, People … ]

[-0.08759557455778122, -0.04312118515372276, -0.08494572341442108, 0.024585919454693794, -0.05785191431641579, -0.02659076638519764, 0.04704275727272034, -0.03940117731690407, 0.005195754114538431, -0.018994472920894623, -0.030896589159965515, -0.02599106915295124, -0.029802896082401276, -0.009517285041511059, -0.03624524921178818, 0.0029738633893430233, -0.04270448908209801, -0.0890769511461258, -0.04064304754137993, 0.01777574978768255, 0.0910411849617958, 0.05333533510565758, -0.07692492008209229, 0.08628936856985092, -0.04232612252354126, -0.007681592833250761, 0.041417201183033, -0.030358949676156044, 0.05717118829488754, 0.0396726056933403, -0.09482061862945557, 0.05382954701781273, -0.01618970558047294, 0.0013550696894526482, 0.004251557867974043, -0.10439810156822205, 0.01734139770269394, 0.08733568340539932, -0.02014184184372425, 0.0690529346466645, -0.05219369381661835, -0.00837920513004646, 0.05078992992639541, -0.0521097406744957, 0.02524719014763832, -0.09064795076847076, -0.0160515438765287, -0.08548879623413086, 0.0957952290773391, 0.0722256377339363, -0.0174766387790441, -0.0711949095129966, 0.0431281477212905, 0.00651295715942978, 0.04662078991532326, 0.05369564145803451, 0.0017072528135031462, 0.0468018501996994, 0.03436211124062538, -0.0556086413562297, 0.04996718838102921, 0.0951215624809265, 0.0048730792477726, -0.0013316450640559196, 0.063605397939682, 0.0242865532636642, 0.00976203475147485, -0.0457112602889537, 0.03228874504566192, -0.0671159550547599, 0.0317926295101642, 0.0003801546408794820, 0.013781636022031307, 0.04971658810973167, -0.04940832778811455, 0.04659050330519676, 0.08707352727651596, -0.10198234766721725, 0.0012964112684130669, 0.01982604712247848, 0.02079876139760017, 0.03980609402060509, -0.0161056233915911, 0.0988031476730687, 0.03530235588550567, -0.03354038670659065, -0.060332611203193665, 0.00999201741069555, -0.0792296230792999, -0.0467287562787532, 0.029249077662825584, 0.00738610792905092, 0.01445483043789863, 0.0367921367287635, 0.0225952658802270, -0.0754463001089096, 0.0375224202871322, -0.0429877676069736, 0.042935766279697, -0.04286639392375946, 0.05493065342307091, -0.010759958997368813, -0.02642848715186119, 0.059368394315242, 0.01073555648326873, -0.0221065618097782, 0.00509836338460445, 0.028178259730339, 0.06781460344791412, 0.01152470614761114, 0.04529837146401405, -0.10414177924394608, -0.063334994018078, -0.02536913566291332, 0.01380544807619076, 0.08876033872365952, -0.0323718637228012, 0.08923118561506271, 0.03578897193074226, -0.073584914207458, -0.11431705206632614, 0.00416507897898554, 0.0398943759500980, 0.01169857848435640, 0.00912736915019371, -0.007194945588707924, -0.06279811263084412, -0.01203678268939256, 0.040613092482089996, -0.07241667807102203, 0.0680359750986099, 0.0318059809505939, 0.020111583173274994, 0.0541915632843366, -0.06192755699157715, 0.0300670601427555, -0.015478908843322943, -0.05181756243109703, -0.02790451049408687, 0.01615394838154316, 0.00118022342224981, -0.052437208592891, -0.0580908432602882, 0.02569353394210338, -0.0445638187229633, -0.05235142633318901, 0.0028490237891674, 0.0251784324645996, 0.07427450269460678, -0.0044151491448283195, 0.0670843124389648, 0.0102450344711542, -0.0372107475960174, -0.05194145068526268, 0.01115098316222429, -0.0136907296255230, -0.0044749984517693, 0.0840763598680496, 0.0184206049889326, 0.0693516209721565, -0.0708598941564, -0.04950730875134468, -0.12054812163114548, -0.0061439946293830, 0.018353456631302834, 0.10588039457798004, 0.04850513488054275, -0.0141451349481940, -0.0234359055757522, -0.0077395206317305565, -0.0335096120834350, 0.02147112786769867, 0.0766637176275253, 0.06256308406591415, 0.01227859035134315, -0.0174938272684812, -0.0258525852113962, 0.0285028610378503, -0.0409292578697204, 0.13894943892955, -0.00052762113045901, 0.019603941589593887, 0.04787249863147335, -0.0247964318841695, -0.0327886901795864, -0.0329894274473190, -0.0749734565651654, 0.0178692955253729, 0.01435997430235147, 0.03260161727666855, -0.0589858926832675, 0.04998610913753509, -0.016536226496100426, -0.0198389366269111, -0.0303456187248229, -0.04567880555987358, -0.0463615022599697, -0.07739298790693283, 0.05453772097826004, -0.015122903152801, -0.0558234192430973, 0.02098872698843479, 0.038055658340454, 0.0410100519657135, 0.06500554829835892, -0.0144112687557937, 0.063243038952506, 0.12235700339078903, -0.00313800736330449, 0.041500613093376, -0.0546980649232864, 0.028155071660876274, -0.02541017346084118, 0.0236208252608776, -0.055749375373125076, 0.034880928695201874, 0.02925102412700653, -0.04492440819740295, -0.0579662285745143, 0.16091801226139069, -0.01372791174799203, -0.0028238031081855297, 0.00505891163359880, 0.0093303751204375, -0.040790800005197525, 0.012631899677217007, -0.0349898150523562, -0.09075972472333562, -0.0406805910170078, 0.0564486208963394, 0.0793653130531311, -0.10637512058030131, 0.024917140916291237, 0.076641976812133, 0.08899114280939102, -0.07655515521764755, 0.0186889953911304, -0.0366524159908294, -0.0071749486960470, -0.0267649386078119, 0.00681320577859878, 0.034349352121353, 0.06390520185232162, 0.00547507125884294, 0.00645425589755177, 0.0124187394976615, -0.0487505868077281, -0.025715982541441917, -0.0013290401548147202, -0.0036538743879646063, -0.034931264817714, 0.07465724647045135, 0.04718988761305809, -0.027499066665768623, 0.011664669029414654, 0.020739786326885223, -0.00119996943976730, 0.0228341985493898, -0.03957423940300941, 0.02761940658092498, 0.01381973270326852, 0.00996309332545972, 0.0177084263414144, -0.0130031989336013794, -0.0075829680180905418, -0.002554208738729357, 0.0359633117914199, 0.0165165513753891, -0.0531400032341480, -0.031333882316165047, 0.07702472805976868, -0.064059898257255, -0.051130335777997, 0.1112219467759132, -0.043169301003217, -0.035172719508409, -0.06082572788000107, -0.057218879461288, -0.0494021698832511, -0.04693160206079483, -0.04560410603880882, -0.0639628246426582, -0.08668574690818787, -0.0083593414241075, 0.08004070073366165, 0.0208362713456153, 0.03663250431418419, 0.0123007167130708, -0.04245587065815925, 0.006077070720493793, 0.003692914033308625, 0.04990659188479185, -0.00278487685136497, -0.03410743912340234, -0.0587813742458203, 0.0523184984924090, 0.01800959184768515, 0.05601579323410988, 0.06174979358911514, 0.00958694703876972, 0.10309688746929169, -0.0013724834425374866, -0.0373494103550910, -0.06568935513496399, -0.0019719060510396957, -0.0252291243523359, 0.00681962072849273, -0.08099189400672913, 0.14139338552141189, 0.025998227298259735, 0.0444508530199527, 0.0934427455067634, -0.0222721491008996, -0.0194158349186182, 0.0552913956344127, -0.03372661769390106, 0.07417678833007812, 0.00704198794632193, -0.053876567631959915, 0.10055916011333852, -0.1129388460369385, -0.033761572092771, 0.06026617065070256, 0.03926900401714, 0.035960644483566284, 0.008666090667247772, 0.027477947995066643, 0.03710789978504181, -0.02351015247404575, 0.035328868776559, -0.0027865190058946, -0.0218553729355335, -0.02822303213179111, 0.04290595650672912, -0.02157396450638771, -0.029180855712767, 0.03941932320594787, -0.0634125620126724, -0.0498380213975906, 0.0411033742129802, -0.0845982655882835, -0.0303522776812315, 0.03796708956360817, 0.1756239682435989, 0.0512967295944690, 0.005838119424879551, -0.05645650997757912, 0.02716675214494333, -0.02447921037673950, -0.05037922412157059, 0.0492484867572784, 0.0166201461110653877, -0.04153639078347607, 0.0292133595794439, 0.040004440456628799, 0.01275655440986156, -0.007722984999418259, -0.03866214305162, 0.01763699576256594, -0.045994330195977271, 0.05585790053009987, -0.008977846242487, -0.07755934447050095, 0.027975516393780, 0.0591881871223449, 0.00961016118526458, 0.09371864050626755, 0.006441071629524231, -0.01596050709486007, 0.10366759449243546, 0.0083675812929868, 0.05804479494690895, 0.04077683761715889, -0.003010124666616320, -0.0076530412770807, 0.12586210668087006, 0.02910740114748478, 0.0408298559486866, -0.08000585434368200, -0.11633670330047607, -0.10950554940760, 0.06872736662626266, -0.0156934037804026, 0.04560477659106254, 0.04468008130788803, 0.0753185658617820, 0.030426912200709771, -0.0503707230091095, 0.00918269529938697, -0.0668076649308204, 0.032497283071279526, -0.038834039121866226, 0.02837871201336383, -0.06321480125188828, -0.00538553716614842, 0.1105560287833213, -0.131211772561073, -0.08003625273704529, -0.00334779825061559, -0.08428100496530533, 0.0205511823965469, -0.02725204080343246, 0.04237489402294159, 0.03843776136636734, -0.01511218305677175, 0.01783227361738681, 0.0724357217550277, 0.03400840610200522, 0.027588229295635223, -0.0224411990493536, -0.0291012115776538, 0.00895093847066164, -0.0130476402118020, 0.0546080631870964, 0.0017075719079002738, -0.03108214400708675, -0.06595993787050247, -0.0136698111891746, -0.01947331987321376, -0.0751525312662124, 0.0577381476759910, 0.00572156766429543, 0.0236212089657783, 0.040137402720451, -0.0762732177997272, -0.04665115848183632, -0.04337655752897262, 0.06120294332504272, 0.00448986794799566, -0.02042182162404060, -0.04481882974505424, -0.003844935214146971, 0.056634843349456, 0.0371888615190982, -0.02901551127433777, -0.0371345877647399, -0.05943121761083603, -0.06698095798492432, 0.0127306943759322, -0.02154143340885639, -0.01820573396980762, 0.01443230827525806, -0.07051131874322891, 0.11712339520454407, 0.0007166709401644766, -0.06536946445703506, -0.0153149710968136, 0.03457680717110634, 0.06144138798117637, 0.0247887931764125, 0.01627664345135, -0.06662845611572266, -0.015765206888318062, -0.003749340074136853, -0.00494204461574554, -0.04201525449752807, -0.07548461109399796, -0.03183511644601822, 0.0612889826297, 0.06246870756149292, -0.00130781123880296, -0.0632513538002967, -0.12265635281801224, 0.02072513103485107, 0.0280458368360996, 0.01681022159755, -0.0145926643162965, 0.043280523270368576, -0.

5. **Test the effectiveness of the document representation through document clustering and classification**

# Dataset: <Reuters>

# Dataset: <Reuters>

**Total Number of Documents: 612,374 (2006. 09. 01 ~ 2015. 06. 06)**

- Random sample of 10,000 documents from 5 distinctive categories (Sports, Market, Politics, Business, World, Health)

| Categories | Total Number of Documents | Number of Sampled Document |
|---|---|---|
| Entertainment | 25,764 | - |
| **Sports** | **49,883** | **10,000** |
| Technology | 26,899 | - |
| **Market** | **189,399** | **10,000** |
| Oddly Enough | 5,864 | - |
| **Politics** | **42,319** | **10,000** |
| **Business** | **96,611** | **10,000** |
| Art | 4,792 | - |
| **World** | **138,852** | **10,000** |
| **Health** | **31,991** | **10,000** |

# Word2Vec & Doc2Vec Training

- During training, words that occur less than 20 times are discarded for stable result
  - Number of unique tokens: 67,390
- **Hyperparameters** to consider while training word2vec / doc2vec models:
  - How many epochs to iterate over the documents?
  - Averaging or concatenating the vectors in the hidden nodes?
  - Number of nodes in hidden layers (**dimension of embedding vectors**)
  - Number of windows (**number of words to use as contexts**)
- No universal hyperparameter settings exist as they are highly dependent on the characteristics and the amounts of the training corpus
- Two types of evaluation methods for word2vec methods:
  1. Extrinsic Evaluation
     - Since the embedded vectors are used as ingredients for building more complex task-specific language model(usually as a pre-training step), evaluation on actual real task
     - Take a long time to compute accuracy
     - Unclear if the subsystem is the problem or its interaction or other subsystems are the problems
  2. Intrinsic Evaluation
     - Evaluation on a specific/intermediate subtask
     - Fast to compute
     - Not helpful unless correlation to real task is established

# Word2Vec & Doc2Vec Training

- Trained 100 different Word2Vec models
  - Dimension: 200 ~ 800
  - Context Window: 6 ~ 10
- Chose the model with the best accuracy in intrinsic evaluation criteria provided by Tomas Mikolov
  - A list of four words with specific relationship is given (19,558 analogies)
    - Capital-Country, Opposing words, Nationalities
  - Given only three words out of four words, test if the model can produce a correct answer
- Chosen hyperparameters:
  - Dimension: 550 & Window Size: 9

```
 1    : capital-common-countries
 2   Athens Greece Baghdad Iraq
 3   Athens Greece Bangkok Thailand
 4   Athens Greece Beijing China
 5   Athens Greece Berlin Germany
 6   Athens Greece Bern Switzerland
 7   Athens Greece Cairo Egypt
 8   Athens Greece Canberra Australia
 9   Athens Greece Hanoi Vietnam
10   Athens Greece Havana Cuba
11   Athens Greece Helsinki Finland
12   Athens Greece Islamabad Pakistan
13   Athens Greece Kabul Afghanistan
14   Athens Greece London England
```

```
find finds generate generates
find finds go goes
find finds implement implements
find finds increase increases
find finds listen listens
find finds play plays
find finds predict predicts
find finds provide provides
find finds say says
find finds scream screams
find finds search searches
find finds see sees
find finds shuffle shuffles
```

# Issues with High Dimensional Clustering

- As the dimension of vectors (data points) grow, normal clustering method doesn't work due to following reasons:

1. **Distance metric becomes useless (no sense of proximity)**
   - Under various data distribution and distance function, ratio of distances of the nearest and farthest neighbors to a given target in high dimensional space is almost 1.
   - Meaningfulness of $L_k$ norm worsens faster with increasing dimensionality for higher values of k

2. **Problems associated with local feature relevance or local feature correlation**
   - Presence of irrelevant features or of correlations among subsets of features heavily influences the appearance of clusters in the full dimensional space
   - Dimension reduction cannot be applied as it only considers one subspace of the original data space in which the clustering can be performed

# 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

**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.

# Choosing K

- Concept of inertia was used for selecting K
  - Inertia: Within-cluster difference with the centroid

$$\sum_{i=0}^{n} \min_{\mu_j \in C} (||x_j - \mu_i||^2)$$

- Cluster the embedded word vectors for various values of K and select the one with the lowest average inertia per cluster
- For each experiment with designated k value, best result out of 300 different trials (different initial points) was compared
- **K = 120 selected**
- Document clustering also used same algorithm except the value of K was fixed at 6

# Clustering Result

- **Accuracy**
  - Perhaps need another representation method for comparison instead of random selection baseline

### Doc2Vec vs. Word2Vec Accuracy

# Cluster Labels

- Name of companies

| | |
|---|---|
| delphi's | 2 |
| navistar's | 2 |
| telstra | 2 |
| agrium | 2 |
| elan's | 2 |
| terra's | 2 |
| tiscali | 2 |
| nonvoting | 2 |
| break-up | 2 |
| avon's | 2 |
| genzyme's | 2 |
| mosaic's | 2 |
| discovery's | 2 |
| wynn | 2 |
| elan | 2 |
| firstgroup | 2 |
| sprint. | 2 |
| drahi | 2 |
| sterlite | 2 |
| unitymedia | 2 |

- People

| | |
|---|---|
| kids | 4 |
| olds | 4 |
| mentors | 4 |
| academic | 4 |
| academia | 4 |
| cooks | 4 |
| enroll | 4 |
| professors | 4 |
| rehab | 4 |
| motivational | 4 |
| workplace | 4 |
| disabilities | 4 |
| telehealth | 4 |
| courses | 4 |
| navigators | 4 |
| underemployme | 4 |
| neediest | 4 |
| educational | 4 |
| bilingual | 4 |
| vaccinations | 4 |
| physicians | 4 |
| autistic | 4 |
| undergraduate | 4 |
| nurses | 4 |

- Negative Words

| | |
|---|---|
| belligerence | 9 |
| brainwashed | 9 |
| unjust | 9 |
| sickening | 9 |
| committing | 9 |
| oaths | 9 |
| tyranny | 9 |
| obscenity | 9 |
| gypsies | 9 |
| inhuman | 9 |
| rapes | 9 |
| intent | 9 |
| betraying | 9 |
| revulsion | 9 |
| neglect | 9 |
| punishable | 9 |
| instigator | 9 |
| deplore | 9 |
| mindless | 9 |
| abuses | 9 |

# Cluster Labels

- Name of countries

| | |
|---|---|
| bratislava | 73 |
| wales | 73 |
| czech | 73 |
| energy-rich | 73 |
| catalonia | 73 |
| angola. | 73 |
| overlord | 73 |
| uae. | 73 |
| singapore | 73 |
| thailand | 73 |
| lesotho | 73 |
| lanka | 73 |
| jordan. | 73 |
| fiji | 73 |
| malta. | 73 |
| indonesia. | 73 |
| andorra | 73 |

- Numbers

| | |
|---|---|
| 272 | 14 |
| 275 | 14 |
| excl | 14 |
| 2025s | 14 |
| 393 | 14 |
| 395 | 14 |
| 398 | 14 |
| 29nov | 14 |
| 270 | 14 |
| 273 | 14 |
| 279 | 14 |
| 305 | 14 |
| 796 | 14 |
| 795 | 14 |
| 790 | 14 |
| 830 | 14 |
| 911 | 14 |
| 199 | 14 |
| 198 | 14 |
| 195 | 14 |
| 194 | 14 |

- Food

| | |
|---|---|
| cooking | 18 |
| intake | 18 |
| concoction | 18 |
| smithfield | 18 |
| kft | 18 |
| wholesome | 18 |
| rice | 18 |
| fruits | 18 |
| lunches | 18 |
| breakfast | 18 |
| h.j. | 18 |
| one-cup | 18 |
| hallucinogenic | 18 |
| shoots | 18 |
| mango | 18 |
| salads | 18 |
| buffet | 18 |
| tortilla | 18 |
| steamed | 18 |

# Final Experiment Setup

- Total number of documents used: 204,000

| Categories | Total Number of Documents | Number of Selected Document |
|:---:|:---:|:---:|
| **Entertainment** | 25,764 | 25,500 |
| **Sports** | 49,883 | 25,500 |
| **Technology** | 26,899 | 25,500 |
| **Market** | 189,399 | 25,500 |
| **Oddly Enough** | 5,864 | - |
| **Politics** | 42,319 | 25,500 |
| **Business** | 96,611 | 25,500 |
| **Art** | 4,792 | - |
| **World** | 138,852 | 25,500 |
| **Health** | 31,991 | 25,500 |

# To Do & Issues

1. Need for another representation method for comparison as a baseline
   - Possible candidates: averaging word embedding vectors, tf-idf

2. Need to compare the accuracy of clustering & classification task given different values of K and the number of hidden nodes

**Accuracy**

**1. Doc2Vec**

**2. Word2Vec Clustering (Best K)**

**3. Word2Vec Clustering (Fixed K)**

**4. Word2Vec Averaging**

**Number of hidden nodes**

3. Need to come up with a method that can automatically create labels for each word2vec clusters
   - Probably based on hypernym found by Wordnet (or using the counts)
   - But the issues still remain with pronouns and stemming
   - More profound background research is needed

4. Find examples of misclassifications in Doc2Vec representation that can be explained by word2vec clusters
   - Qualitatively substantiating the explanatory power of the suggest method

5. Finish the experiment for the final test settings and submit a preliminary results for the conference in November

# 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.

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..

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