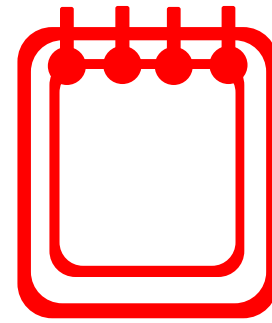


Statistical Computing and Simulation -Final Report



Professor: 翁久幸

Student: 統研一 董承

Procedure

- EM & K-means
- NIPS Clustering

EM & K-means

EM Algorithm

1. 初始化參數

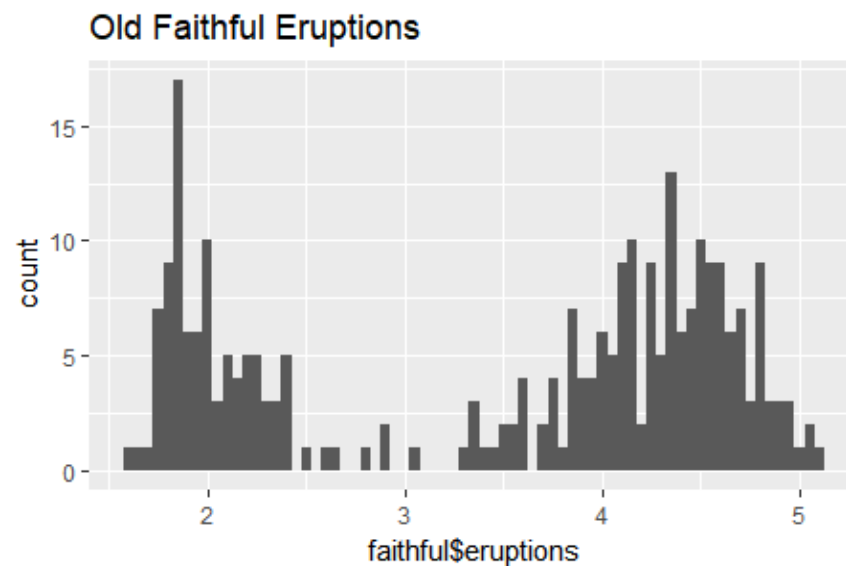
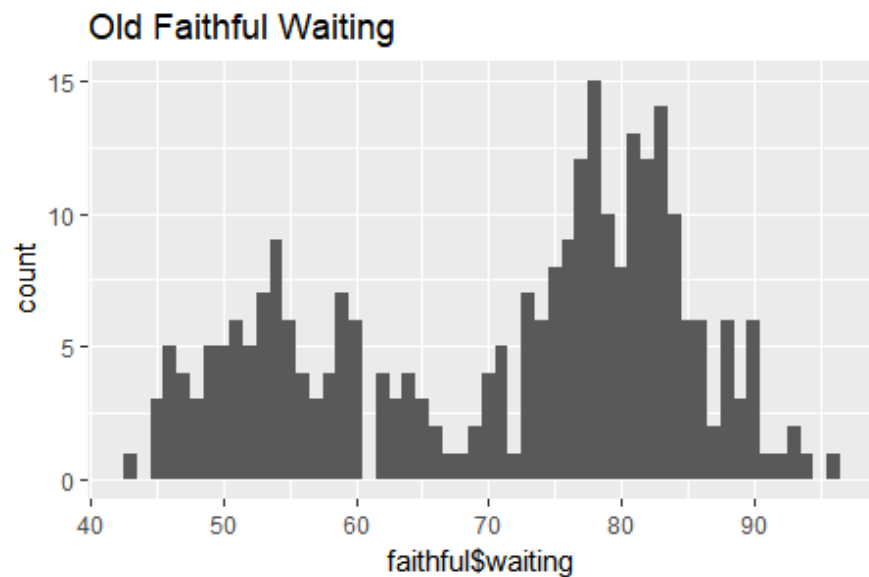
2. 重複直到收斂：

E步驟：估計未知參數的期望值，給出當前的參數估計。

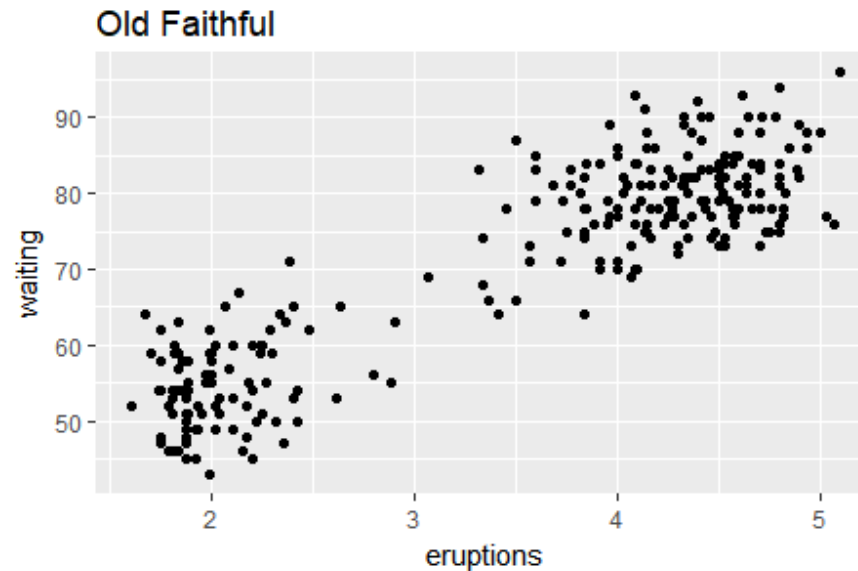
M步驟：重新估計分布參數，以使得最大概似函數最大，得到未知變量的期望估計。

Old Faithful Dataset

- 利用美國黃石國家公園的老忠實噴泉作為比較K-means & EM 資料，包含兩個變數，分別是eruptions(噴發持續時間)、waiting(兩次噴發之間隔時間)，這裡利用視覺化探索資料得到資料分布圖形：



Old Faithful Dataset



- 可以看出這筆資料大致有兩個群，分別是eruptions(噴發持續時間)及waiting(噴發間隔時間)“短”，和eruptions(噴發持續時間)及waiting(噴發間隔時間)“長”。

Expectation– maximization

- 作為missing data的指標變數為

$$Y_i = \begin{cases} 1 & W_i \text{ belongs to distribution of shorter waiting times} \\ 0 & W_i \text{ belongs to distribution of longer waiting times} \end{cases}$$

- Y_i 服從二項分配，參數為 p ，接著，利用EM演算法找到 Y_i 的期望值
- 透過重覆疊代得到 Y_i 的參數 p ，我們可以得到期望值函式，此即為EM演算法的E-step，再利用微分技巧找出每個參數，使得概似函數最大，為EM演算法的M-step
- 而初始值的假設是透過資料探索的圖做為參考，假設如下

$$p^{(0)} = 0.4, \mu_1^{(0)} = 40, \mu_2^{(0)} = 90, \sigma_1^{(0)} = 4, \sigma_2^{(0)} = 4$$

Expectation-maximization

- 分群結果:



Expectation– maximization

```
$params
      Comp.1  Comp.2
means  54.195591 80.353490
variances 4.956753 7.547971
```

```
means  54.194080 80.35394
variances 4.946035 7.53402
```

```
$clusterpobs
  probMembership
1      0.3083102
2      0.6916898
  probMembershipFlex
1      0.3087248
2      0.6912752
```

```
$params
      Comp.1  Comp.2
means  4.2733696 2.0186404
variances 0.4378327 0.2361129
```

```
means  2.0186078 4.2733434
variances 0.2356218 0.4370631
```

```
$clusterpobs
  probMembership
1      0.3484046
2      0.6515954
  probMembershipFlex
1      0.6507353
2      0.3492647
```

- 來自“短”群的waiting與eruptions之平均值分別為54.1956, 4.2733，而資料分別由waiting與eruptions的變數下，屬於“短”群的機率分別為0.3083, 0.3484，反之亦然。

K-means

- 給定K 個clusters,指標變數為

$$r_{i,k} = \begin{cases} 1 & \text{if } x_i \text{ belongs to cluster } k \\ 0 & \text{otherwise} \end{cases}$$

- 由於K means 是屬於 hard clustering ,

也就是每筆資料只有屬於或不屬於k-th 群的選擇 ,

每一點都被分配到一個且是唯一的一個群

K-means

- 首先，隨機取 k 個 seed point，將與第 i 個 seed point 的點距離最近的點分配至群 i ，至所有點分配完為止。這邊所使用的距離測量方式為 Euclidean distance

$$\sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$

- 指標變數可以改寫成：

$$r_{i,k} = \begin{cases} 1 & \text{if } d(\mathbf{x}_i, \boldsymbol{\mu}_k) = \min\{d(\mathbf{x}_i, \boldsymbol{\mu}_1), d(\mathbf{x}_i, \boldsymbol{\mu}_2), \dots, d(\mathbf{x}_i, \boldsymbol{\mu}_K)\} \\ 0 & \text{otherwise} \end{cases}$$

K-means

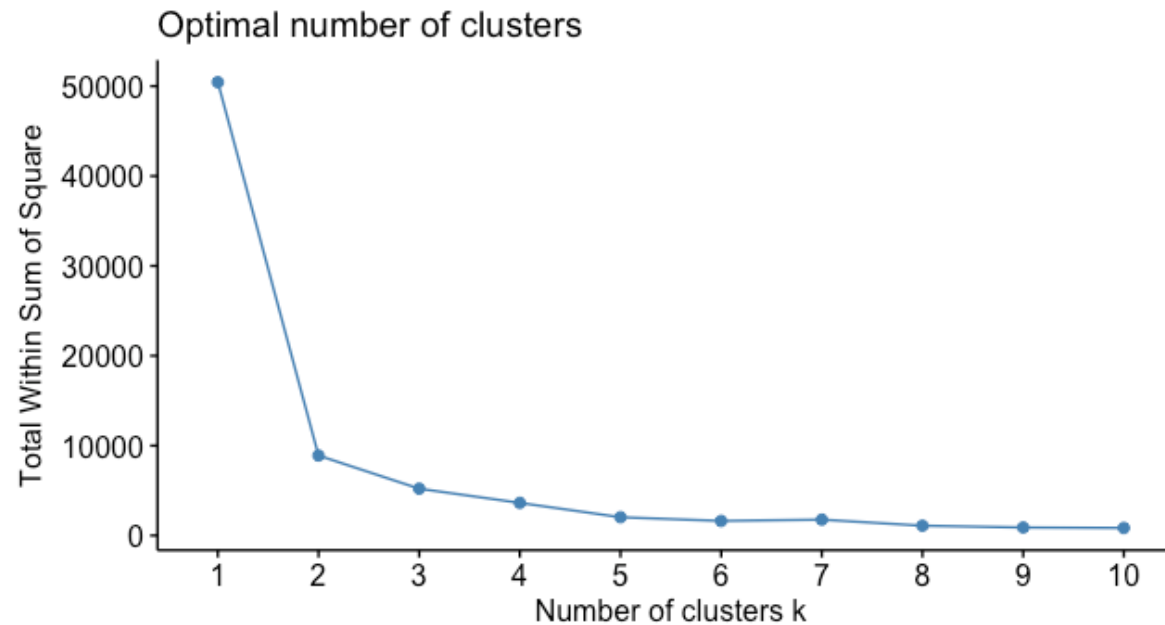
- 計算所形成的群中更新過後的centroids， $k = 1, \dots, K$

$$\mu_k = \frac{1}{\sum_{i=1}^N r_{i,k}} \sum_{i=1}^N r_{i,k} \mathbf{X}_i$$

- 重複上述過程直到每個 centroids 收斂至一特定值。

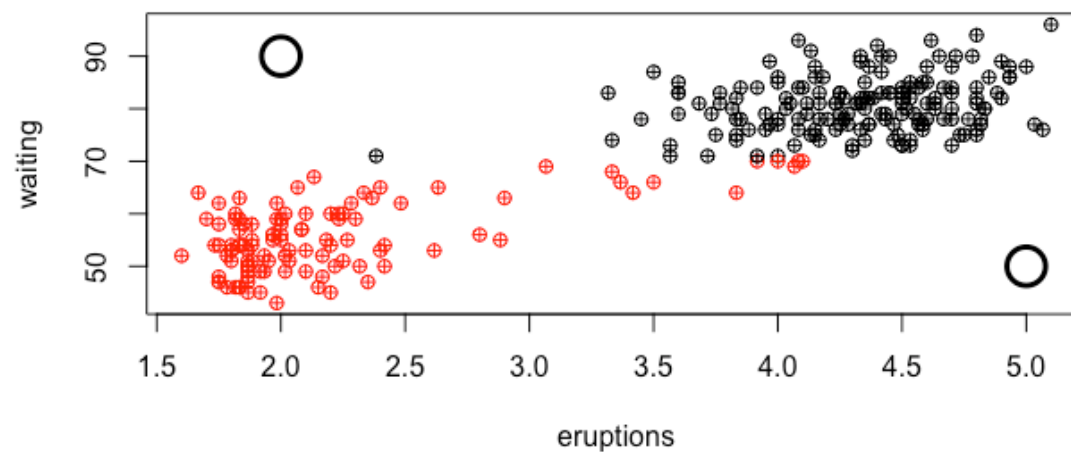
K-means

- 因為K 是執行之前要先給定的值，通常會用組內平方和(within group sum of squares)來決定K 的取值，以下是給定K=1到10 中的組內平方和：



K-means

- 隨著K 值增加，組內平方和會下降，但其中又以 $K=2$ 時，下降幅度最大且與 $K=3$ 時差異不會太大，因此我們選定 $K=2$ ，任取2個點 $(2,90)$, $(5,50)$ 作為起始位置，經過一次迭代，得到：

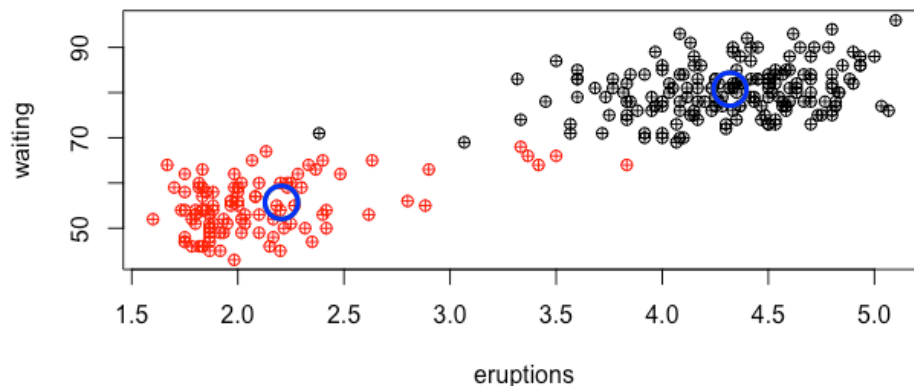


K-means

- 根據此次分群結果再重新計算此兩群內的新中心點：

1. $(2.205607, 55.71028)$ 2 $(4.319255, 80.74545)$

- 計算其他所有點到這兩個中心分別的 Euclidean distance，取距離近的加入該群，第二次迭代後得到的結果如下：



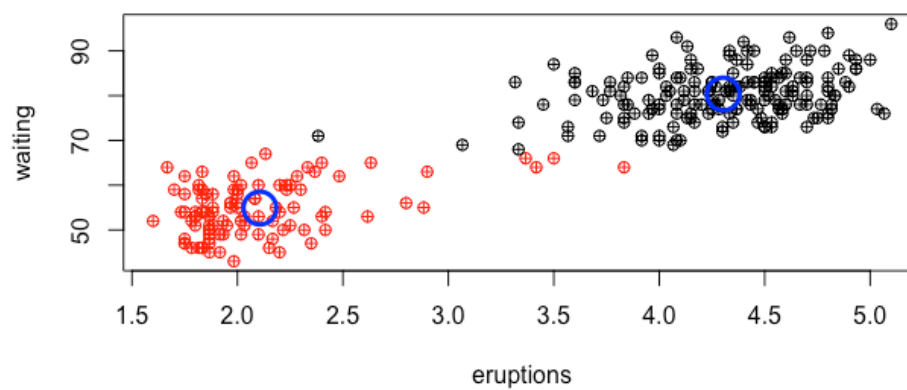
K-means

- 此次分群結果兩群內的新中心點為：

1. (2.106594 , 54.88119)

2. (4.303573 , 80.35673)

- 執行第三次迭代得到以下結果：



K-means

- 由此三次迭代可以發現中心點位置漸漸收斂，使用R 內建 `kmeans()` 執行 `old faithful`(老忠實)噴泉分群，得出以下視覺畫圖型及分群結果：

```
> kmeans.cluster$centers
```

```
eruptions waiting
```

```
1  4.29793 80.28488
```

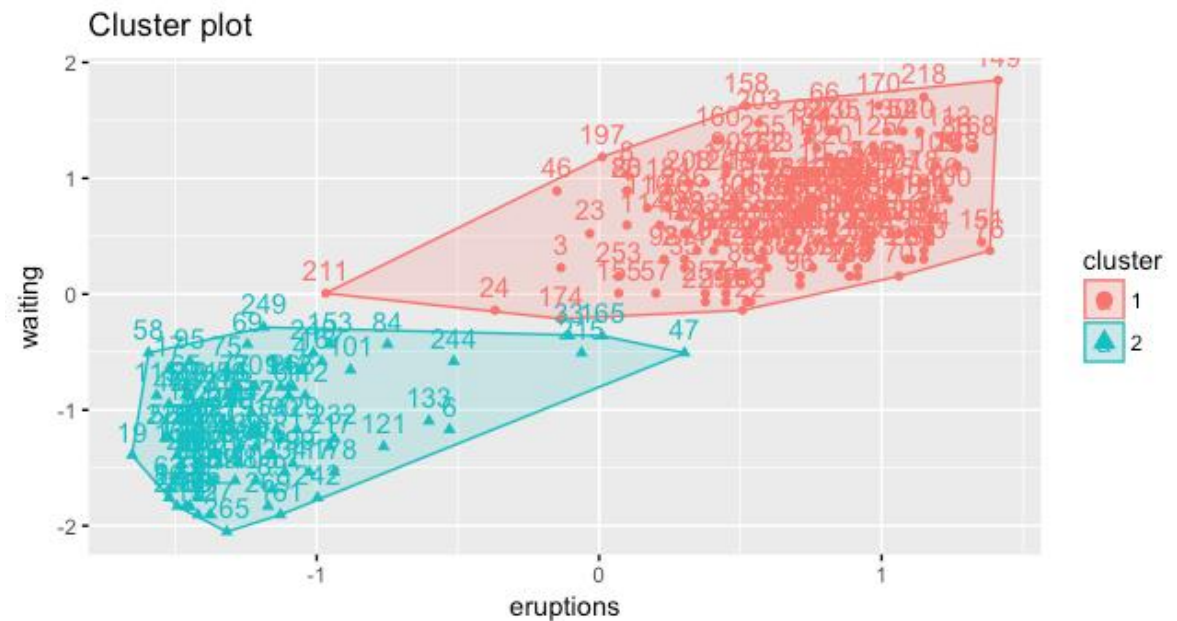
```
2  2.09433 54.75000
```

```
> kmeans.cluster$size
```

```
[1] 172 100
```

```
> kmeans.cluster$withinss
```

```
[1] 5445.591 3456.178
```



Summary

- EM method優點：

1. 可靠地找到局部最優的收斂值
2. 比K-means穩定、準確

- EM method缺點：

1. 對初始值敏感
2. 計算複雜且收斂慢

- K-means 優點：

1. 原理簡單，實現容易
2. 收斂速度快
3. 算法的可解釋度比較強
4. 只有參數k 需要調整

- K means 缺點：

1. k 值不易選取
2. 採用迭代方法，得到的結果只是局部最優
3. 對outliers 較為敏感
4. 群集不能重疊
5. 數據分布情況、群集中心的初始位置都會影響重複次數

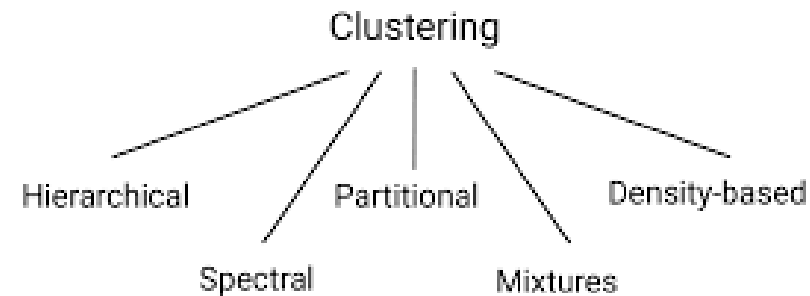
NIPS Clustering

NIPS

- 資料集為NIPS(Conference on Neural Information Processing Systems)神經資訊處理系統研討會的文件資料
- 包含11463個詞以及5812份文件的Term-Document Matrix
- 透過斷詞、移除停用字符，記錄每篇文件中出現詞的頻率，並且只保留那些出現次數高於50次的詞
- 稀疏矩陣(Sparse Matrix)

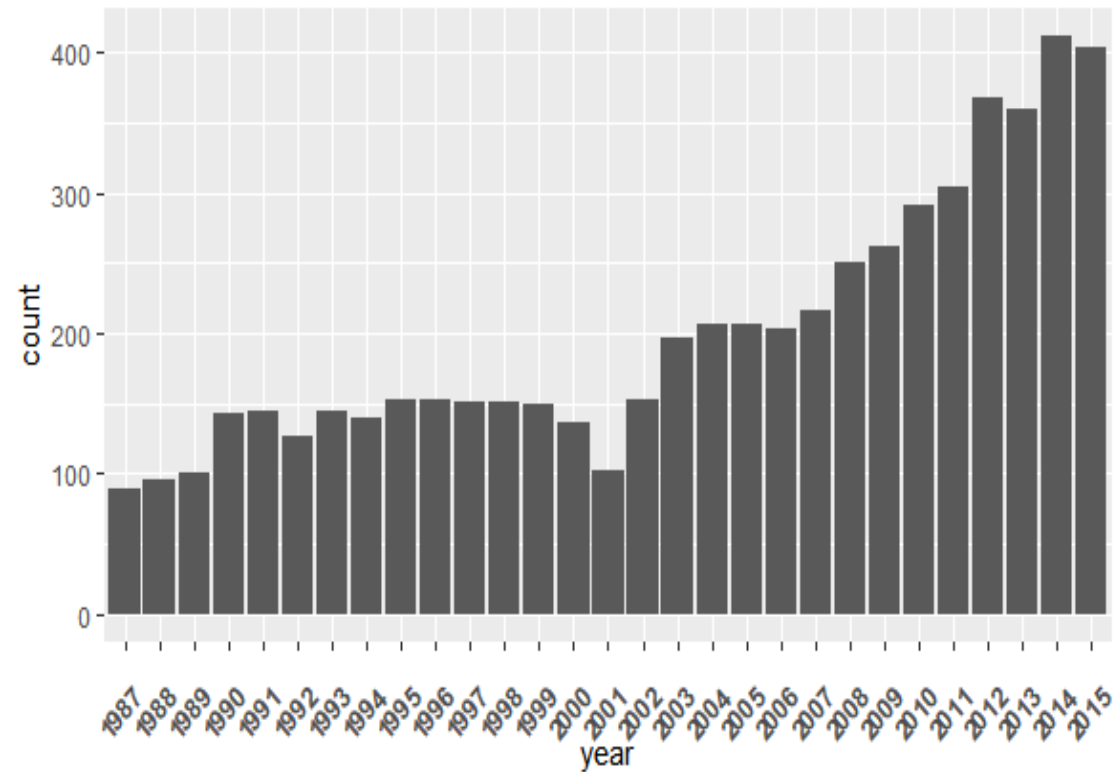
Target

- Cluster the documents by term : Documents with similar key terms are supposed to be in same group.



EDA(Exploratory Data Analysis)

- Number of documents since 1987-2015



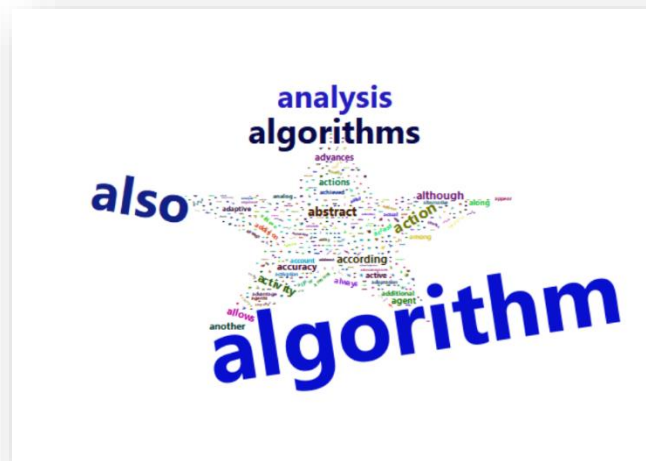
Wordcloud

- R package : wordcloud2

Since 1987-1995



Since 1996-2005



Since 2006-2015



TF-IDF

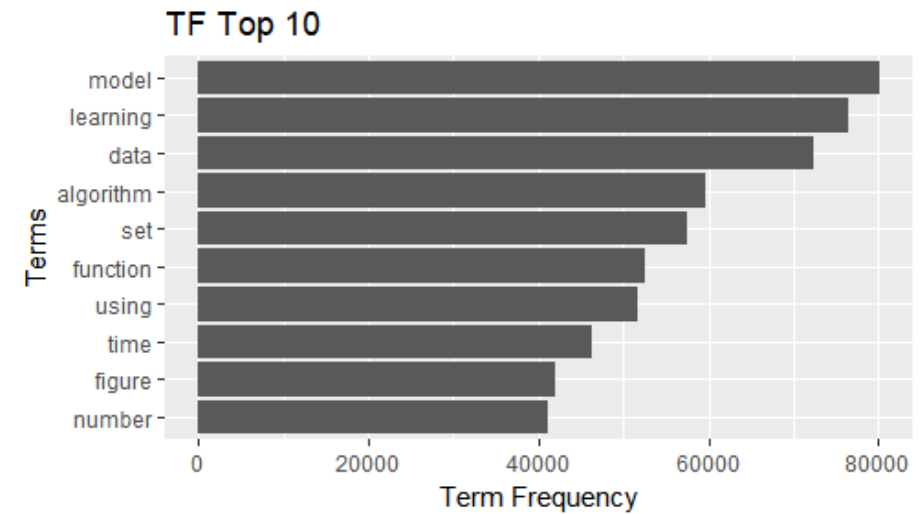
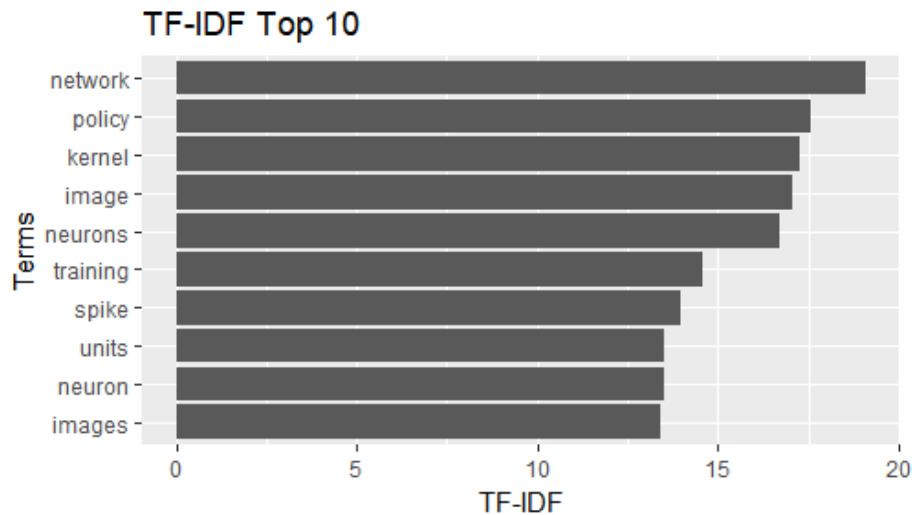
- TF-IDF(Term Frequency–Inverse Document Frequency)
- 詞頻 (Term Frequency , TF) 逆向文件頻率 (Inverse Document Frequency , IDF)

$$TFIDF = TF \times IDF = n_w^d \times \log_2 \left[\frac{N}{N_w} \right]$$

- 其中, n_w^d 表示在文章d中，文字w出現的次數，N表示總共有幾篇文章， N_w 表示有文字w的文章有幾篇。

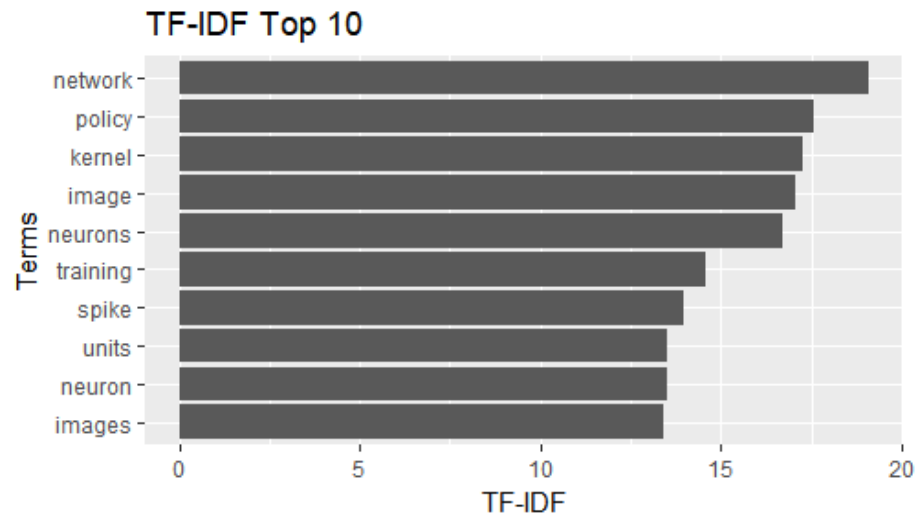
TF-IDF

- R package : tm (text mining)
- By sum the TF-IDF given all documents

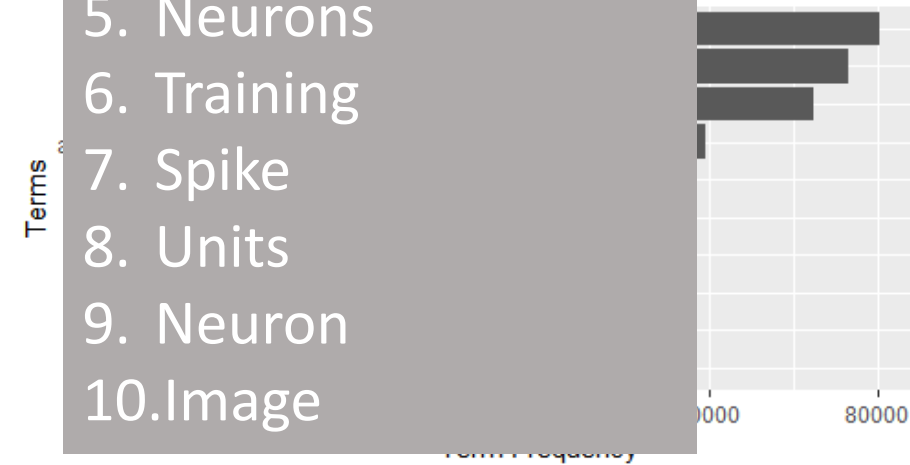


TF-IDF

- R package : tm (text mining)
- By sum the TF-IDF given all documents



1. Network
2. Policy
3. Kernel
4. Image
5. Neurons
6. Training
7. Spike
8. Units
9. Neuron
10. Image



TF-IDF

- See correlated words

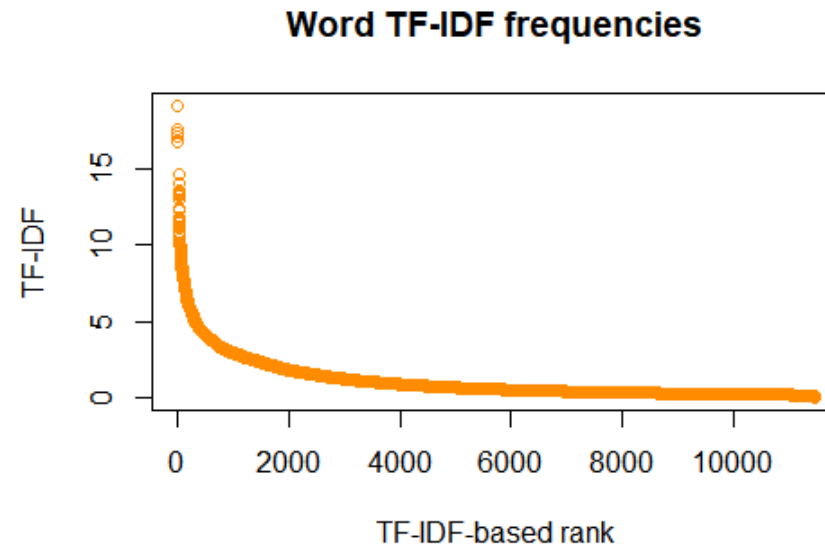
```
> # find correlated terms
> findAssocs(tdm.tfidf, "network", 0.5)
$network
networks    neural
      0.65      0.52
```

```
> findAssocs(tdm.tfidf, "spike", 0.5)
$spike
spikes trains
      0.71    0.67
```

```
> findAssocs(tdm.tfidf, "units", 0.5)
$units
unit hidden
      0.69    0.54
```

Data Manipulation

- Sum the value of TF-IDF given all documents



Data Manipulation

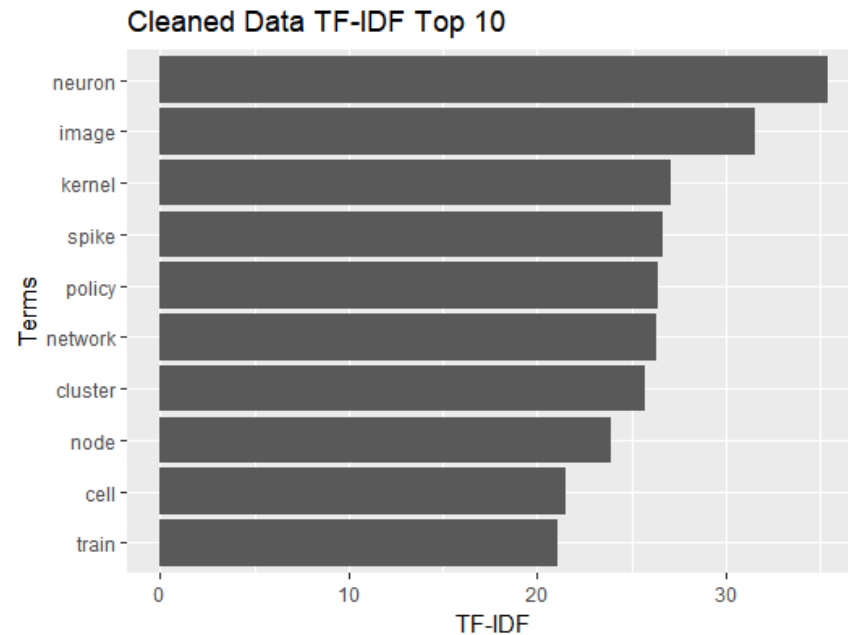
- Stemming & Lemmatisation
- R command : `textstem::lemmatize_words()`
- `c("slowly" , "ran" , "seen" , "eats") %>% lemmatize_words()`

`c("slowly", "run", "see", "eat")` $\left[\right]_{2000 \times 5811}$

- R command : `plyr::ddply` $\left[\right]_{1453 \times 5811}$

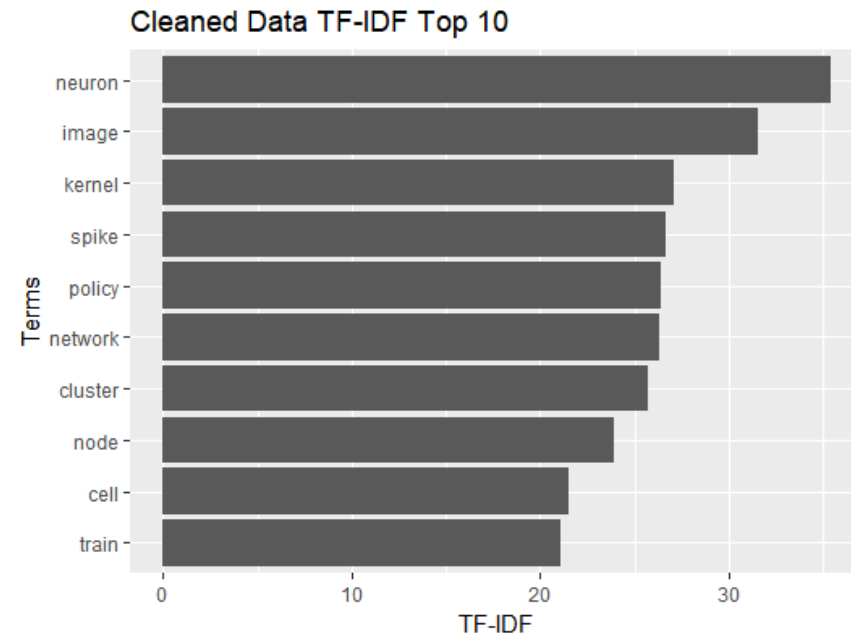
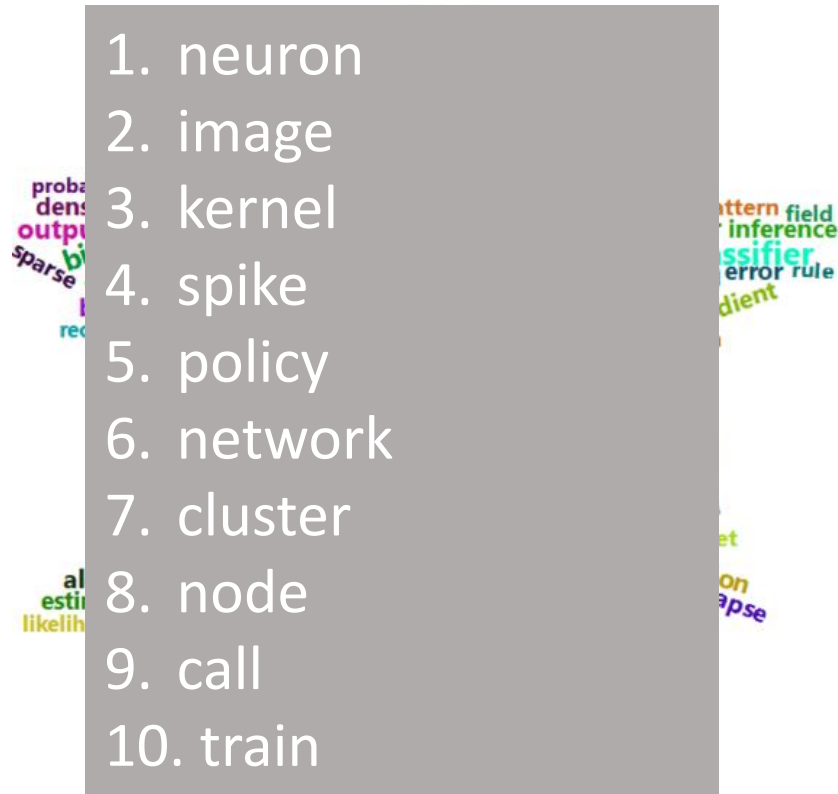
1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2632, 2633, 2634, 2635, 2636, 2637, 2638, 2639, 2640, 2641, 2642, 2643, 2644, 2645, 2646, 2647, 2648, 2649, 2650, 2651, 2652, 2653, 2654, 2655, 2656, 2657, 2658, 2659, 2660, 2661, 2662, 2663, 2664, 2665, 2666, 2667, 2668, 2669, 2670, 2671, 2672, 2673, 2674, 2675, 2676, 2677, 2678, 2679, 26

- Abstract**



Data Manipulation

- WordCloud & Top 10 TF-IDF Terms



Vector Space Model

- 餘弦相似性(Cosine Similarity)是通過測量兩個向量的夾角的餘弦值來度量它們之間的相似性。0度角的餘弦值是1 (兩個向量有相同的指向時)；當兩個向量夾角為90°時，餘弦相似度的值為0；而兩個向量指向完全相反的方向時，餘弦相似度的值為-1。其計算方式如下

$$\text{Similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

[]
5811×5811

Vector Space Model

- Top 20 Documents Related to X1987_1

```
> aa[5792:5811,1:2] # see top 20 documents  
related to X1987_1
```

Docs	Docs	X1987_1	X1987_2
X2003_158		0.4181841	0.3280711
X1992_123		0.4187696	0.2599045
X2002_133		0.4204830	0.2879627
X2005_94		0.4226493	0.3560670
X2015_230		0.4356350	0.5413217
X2000_1		0.4404999	0.3223251
X2000_2		0.4409303	0.3220780
X2006_77		0.4445143	0.4487613
X1997_100		0.4486664	0.5279541
X2014_120		0.4549942	0.2833302
X1993_65		0.4638252	0.4223435
X2007_52		0.4646202	0.4228044
X1995_30		0.4720592	0.3845963
X1999_16		0.4723698	0.3793990
X1992_71		0.4899850	0.3513951
X1993_62		0.4958996	0.3648531
X1987_30		0.5029856	0.3273626
X1992_100		0.5141574	0.4492241
X1990_143		0.5305098	0.4073333
X1987_1		1.0000000	0.3300572

Clustering

• 階層式分群法 (Hierarchical Clustering)

Names	Formula
Euclidean distance	$\ a - b\ _2 = \sqrt{\sum_i (a_i - b_i)^2}$
Squared Euclidean distance	$\ a - b\ _2^2 = \sum_i (a_i - b_i)^2$
Manhattan distance	$\ a - b\ _1 = \sum_i a_i - b_i $
Maximum distance	$\ a - b\ _\infty = \max_i a_i - b_i $
Mahalanobis distance	$\sqrt{(a - b)^\top S^{-1} (a - b)}$ where S is the covariance matrix

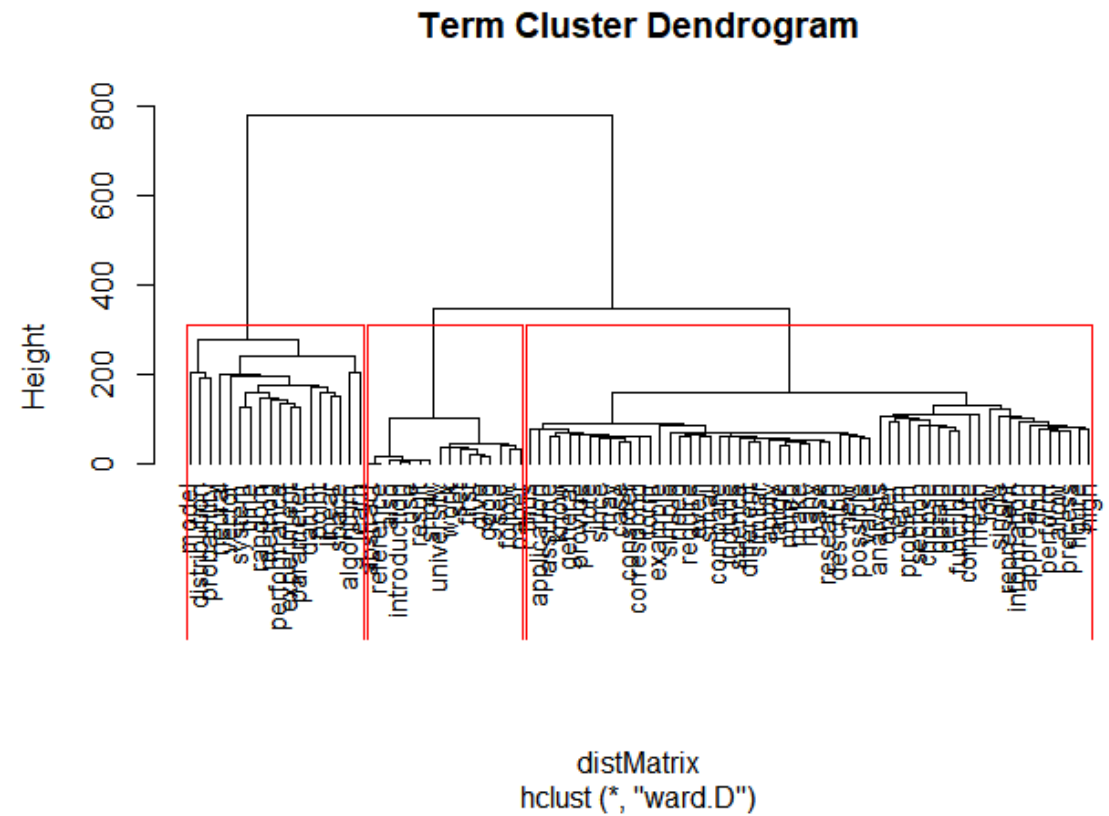
Names	Formula
Maximum or complete-linkage clustering	$\max \{ d(a, b) : a \in A, b \in B \}.$
Minimum or single-linkage clustering	$\min \{ d(a, b) : a \in A, b \in B \}.$
Mean or average linkage clustering, or UPGMA	$\frac{1}{ A B } \sum_{a \in A} \sum_{b \in B} d(a, b).$
Centroid linkage clustering, or UPGMC	$\ c_s - c_t\ $ where c_s and c_t are the centroids of clusters s and t respectively.
Minimum energy clustering	$\frac{2}{nm} \sum_{i,j=1}^{n,m} \ a_i - b_j\ _2 - \frac{1}{n^2} \sum_{i,j=1}^n \ a_i - a_j\ _2 - \frac{1}{m^2} \sum_{i,j=1}^m \ b_i - b_j\ _2$

where d is the chosen metric. Other linkage criteria include:

- The sum of all intra-cluster variance.
- The increase in variance for the cluster being merged (Ward's criterion).^[7]
- The probability that candidate clusters spawn from the same distribution function (V-linkage).
- The product of in-degree and out-degree on a k-nearest-neighbour graph (graph degree linkage)

Clustering

- Hierarchical clustering on Terms

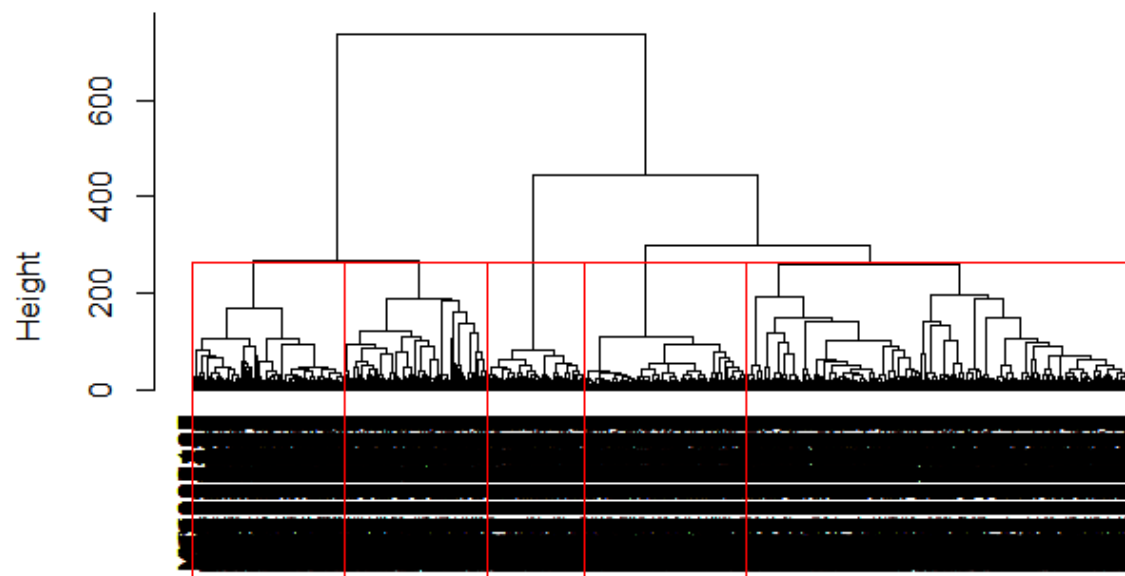


Clustering

- Hierarchical clustering on Documents

```
> bb[which(cut.h.cluster == 1),1] %>% names() # x1987_1
[1] "x1987_1" "x1987_5" "x1987_7" "x1987_9" "x1987_10"
[6] "x1987_11" "x1987_19" "x1987_25" "x1987_27" "x1987_33"
[11] "x1987_34" "x1987_37" "x1987_39" "x1987_41" "x1987_42"
[16] "x1987_43" "x1987_45" "x1987_46" "x1987_49" "x1987_51"
[21] "x1987_55" "x1987_57" "x1987_58" "x1987_63" "x1987_66"
[26] "x1987_68" "x1987_79" "x1987_81" "x1987_88" "x1987_90"
[31] "x1988_15" "x1988_21" "x1988_22" "x1988_23" "x1988_24"
[36] "x1988_25" "x1988_29" "x1988_33" "x1988_35" "x1988_38"
[41] "x1988_39" "x1988_42" "x1988_46" "x1988_51" "x1988_53"
[46] "x1988_61" "x1988_62" "x1988_64" "x1988_65" "x1988_66"
[51] "x1988_73" "x1988_74" "x1988_80" "x1988_81" "x1988_84"
[56] "x1988_90" "x1988_92" "x1988_93" "x1989_9" "x1989_12"
[61] "x1989_13" "x1989_16" "x1989_20" "x1989_23" "x1989_25"
[66] "x1989_30" "x1989_32" "x1989_42" "x1989_46" "x1989_47"
[71] "x1989_57" "x1989_58" "x1989_61" "x1989_62" "x1989_73"
[76] "x1989_77" "x1989_79" "x1989_85" "x1989_86" "x1989_90"
[81] "x1990_1" "x1990_4" "x1990_7" "x1990_8" "x1990_13"
[86] "x1990_32" "x1990_33" "x1990_34" "x1990_40" "x1990_41"
[91] "x1990_42" "x1990_44" "x1990_48" "x1990_56" "x1990_61"
```

Document Cluster Dendrogram



distMatrix
hclust (*, "ward.D")

Summary

- Data manipulation 清理副詞
- NLP 分群的部分可以做語意分析(LSA)
- Effects of Distance Metrics on Document Clustering

ABSTRACT

Effects of Distance Metrics on Document Clustering

by
Rushikesh Veni

Dr. Kazem Taghva, Examination Committee Chair
Professor, Department of Computer Science
University of Nevada, Las Vegas

Document clustering or unsupervised document classification is an automated process of grouping documents with similar content. A

**THANK
YOU**

