# HW 4 – Exploring Common Topics in Health Research News

Shivangi Mundhra | SUID – 842548148

For the three algorithms, kMeans, LDA and BERTopic, I did not sample the data from the given dataset of 10,000 headlines from health research press releases posted on the EurekAlert! Website. All headlines were included in all training all models to explore the common themes of health research news by applying kMeans, LDA, and BERTopic algorithms.

## kMeans

### Number of clusters –
I used k = 10 to begin with. It resulted in 10 clusters of which cluster #2 was the largest with 4,322 instances. The labels are as below. After using the elbow method to find the best k, I received k = 9 because the slope decreased at k = 9. I used this with sBERT and received below labelled clusters.

### Labels –

| cluster # | kMeans (k = 10) | sBERT (k = 9) |
|---|---|---|
| 1 | Coronavirus | Research and study findings |
| 2 | Tobacco | Depression and Mental Health |
| 3 | Unsure | Obesity |
| 4 | Brain | Pregnancy |
| 5 | Obesity | Cancer |
| 6 | Heart Health | Ifectious disease Pandemic |
| 7 | Breast cancer | Cognitive Diseases |
| 8 | Hysteria | Heart Health |
| 9 | Stroke | Tobacco |
| 10 | Heart Health | -- |

There was a lot of overlap in records between the clusters in kMean (10) algorithm but the sBERT model resulted in precise results. The records included in each cluster centered around a specific topic.

## LDA

### Number of clusters –
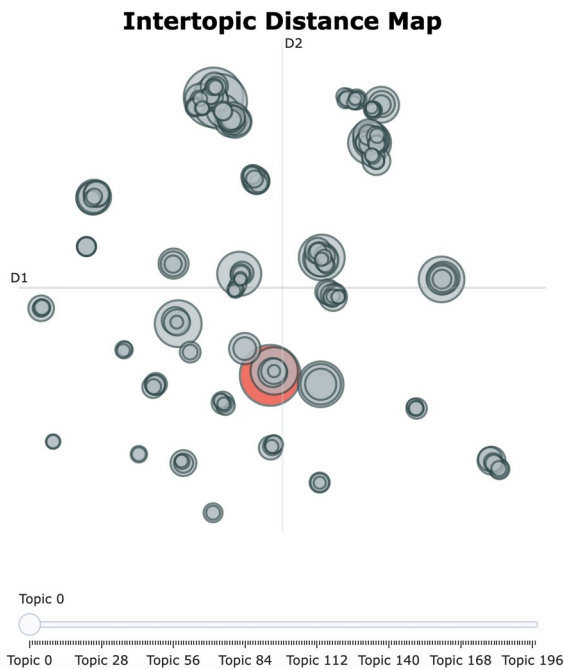I used number of topics = 15.

### Labels –
The below table shows the top 10 words for each label in the SVM model.

| Topic# | Topic |
|---|---|
| 1 | Adult care |
| 2 | Research and findings |
| 3 | Heart Risks |

| 4 | Sickness caused by pollution |
|---|---|
| 5 | Orthopedic ailments |
| 6 | Cancer |
| 7 | STDs |
| 8 | Tests and breakthroughs |
| 9 | Organs |
| 10 | Blood related |
| 11 | Parenthood |
| 12 | Unsure |
| 13 | Covid-19 |
| 14 | Meals |
| 15 | Brain |

## BERTopic

I received 198 topic using BERTopic. Obviously, it is difficult to label all of these. The below chart shows the intertopic Distance Map obtained from BERTopic.



Below is the table of labels.

| topic# | topic | topic# | topic | topic# | topic | topic# | topic |
|---|---|---|---|---|---|---|---|
| 1 | breast | 50 | adhd | 99 | leukemia | 148 | transfusions |
| 2 | dementia | 51 | older | 100 | pulmonary | 149 | sitting |
| 3 | preterm | 52 | ptsd | 101 | expectancy | 150 | exercise |
| 4 | diabetes | 53 | migraine | 102 | sepsis | 151 | valve |
| 5 | tobacco | 54 | pancreatic | 103 | bladder | 152 | violence |
| 6 | gut | 55 | arthritis | 104 | kidney | 153 | american |
| 7 | alcohol | 56 | osteoarthritis | 105 | preterm | 154 | aneurysm |
| 8 | prostate | 57 | epilepsy | 106 | media | 155 | football |
| 9 | heart | 58 | end | 107 | heat | 156 | radiation |
| 10 | opioid | 59 | mental | 108 | cystic | 157 | celiac |
| 11 | eye | 60 | errors | 109 | breast | 158 | toxicity |
| 12 | hiv | 61 | medical | 110 | cervical | 159 | dengue |
| 13 | antibiotic | 62 | trials | 111 | fibrillation | 160 | secondhand |
| 14 | heart | 63 | statins | 112 | zika | 161 | lymphoma |
| 15 | liver | 64 | hot | 113 | arrest | 162 | autoimmune |
| 16 | tumor | 65 | stimulation | 114 | traumatic | 163 | cocaine |
| 17 | vitamin | 66 | apnea | 115 | withdrawal | 164 | incontinence |
| 18 | stroke | 67 | back | 116 | fruit | 165 | tamoxifen |
| 19 | colorectal | 68 | pregnancy | 117 | reconstruction | 166 | alcohol |
| 20 | lung | 69 | preeclampsia | 118 | sleep | 167 | coffee |
| 21 | obesity | 70 | ebola | 119 | spinal | 168 | dismissing |
| 22 | asthma | 71 | smoking | 120 | psoriasis | 169 | appendicitis |
| 23 | hip | 72 | depression | 121 | falls | 170 | hepatitis |
| 24 | cannabis | 73 | fat | 122 | breast | 171 | clock |
| 25 | covid | 74 | hpv | 123 | postpartum | 172 | nsclc |
| 26 | food | 75 | hearing | 124 | exercise | 173 | mindfulness |
| 27 | brain | 76 | gun | 125 | pain | 174 | particulate |
| 28 | flu | 77 | dialysis | 126 | depressed | 175 | violent |
| 29 | autism | 78 | kidney | 127 | thyroid | 176 | smokers |
| 30 | pressure | 79 | gestational | 128 | hands | 177 | hiv |
| 31 | malaria | 80 | childhood | 129 | copd | 178 | burnout |
| 32 | sleep | 81 | clots | 130 | ketamine | 179 | thinners |
| 33 | schizophrenia | 82 | mammogram | 131 | older | 180 | antidepressa |
| 34 | hepatitis | 83 | eating | 132 | salt | 181 | gout |
| 35 | melanoma | 84 | meat | 133 | nurses | 182 | radionuclide |
| 36 | ovarian | 85 | road | 134 | fitness | 183 | knee |
| 37 | bariatric | 86 | obesity | 135 | medicaid | 184 | chocolate |
| 38 | sexual | 87 | esophageal | 136 | stents | 185 | muscular |
| 39 | exercise | 88 | neck | 137 | breastfeedin | 186 | fatigue |
| 40 | kidney | 89 | racial | 138 | walking | 187 | media |
| 41 | tb | 90 | ct | 139 | gender | 188 | transplant |
| 42 | allergy | 91 | testosterone | 140 | financial | 189 | shoulder |
| 43 | parkinson | 92 | readmission | 141 | older | 190 | coronavirus |
| 44 | plastic | 93 | delirium | 142 | work | 191 | pet |
| 45 | sclerosis | 94 | dna | 143 | erectile | 192 | diet |
| 46 | sugar | 95 | omega | 144 | organ | 193 | abuse |
| 47 | suicide | 96 | countries | 145 | anesthesia | 194 | pollution |
| 48 | pollution | 97 | concussion | 146 | parasite | 195 | body |
| 49 | cholesterol | 98 | tooth | 147 | eczema | 196 | rebuilding |
|  |  |  |  |  |  | 197 | obesity |
|  |  |  |  |  |  | 198 | depression |

## Ethic Statement

There are certainly overlaps in many of the clusters in probably all these models. We understand that all these health topics are sensitive and sometimes controversial. The idea is to create a best estimate for research and study purposes and not to slight any article about any disease that we may have missed to cast light on.

## Conclusion

The effectiveness and precision by which BERTopic classifies all these headlines so fast is impressive. None of the other models produce such a fine distinction between headlines of topics. Plus BERTopic also eliminates the manual effort to guesstimate the number of clusters which makes it more desirable.