# DTSA 5510 Unsupervised Learning Final Project
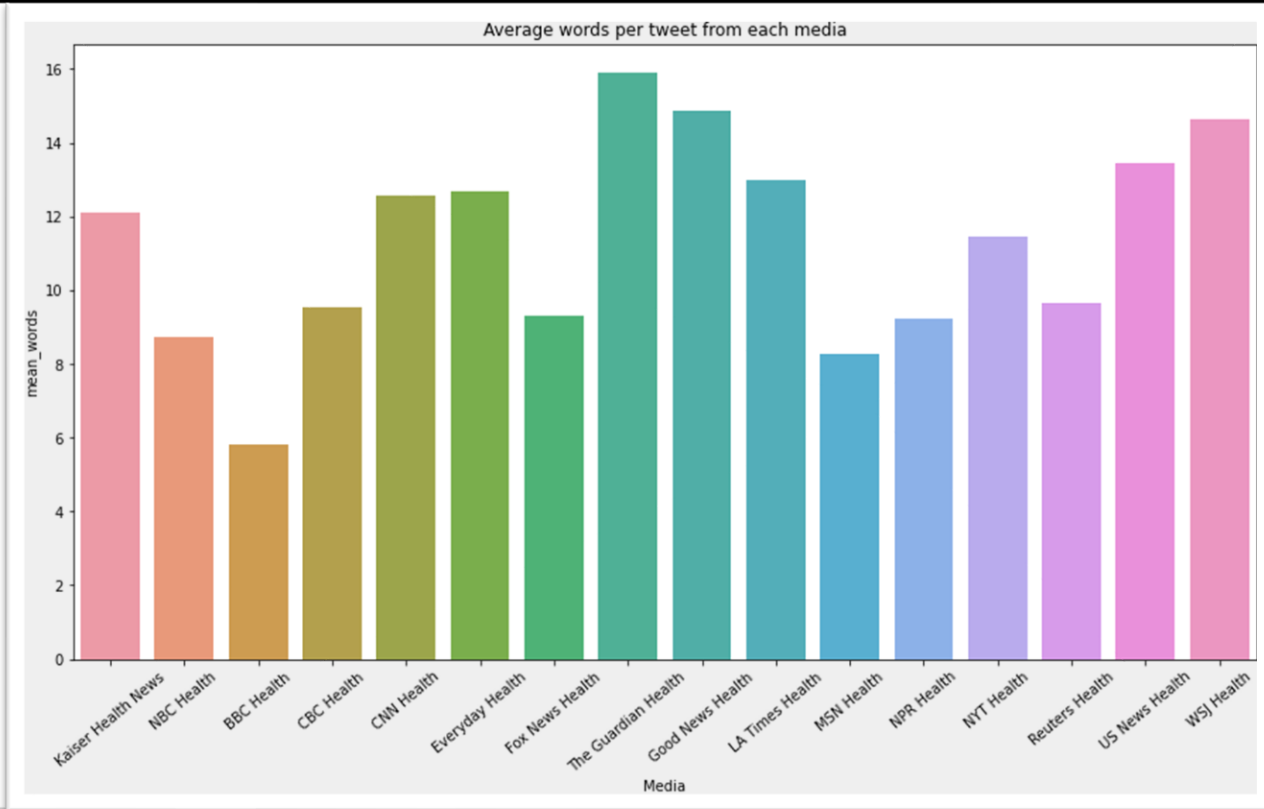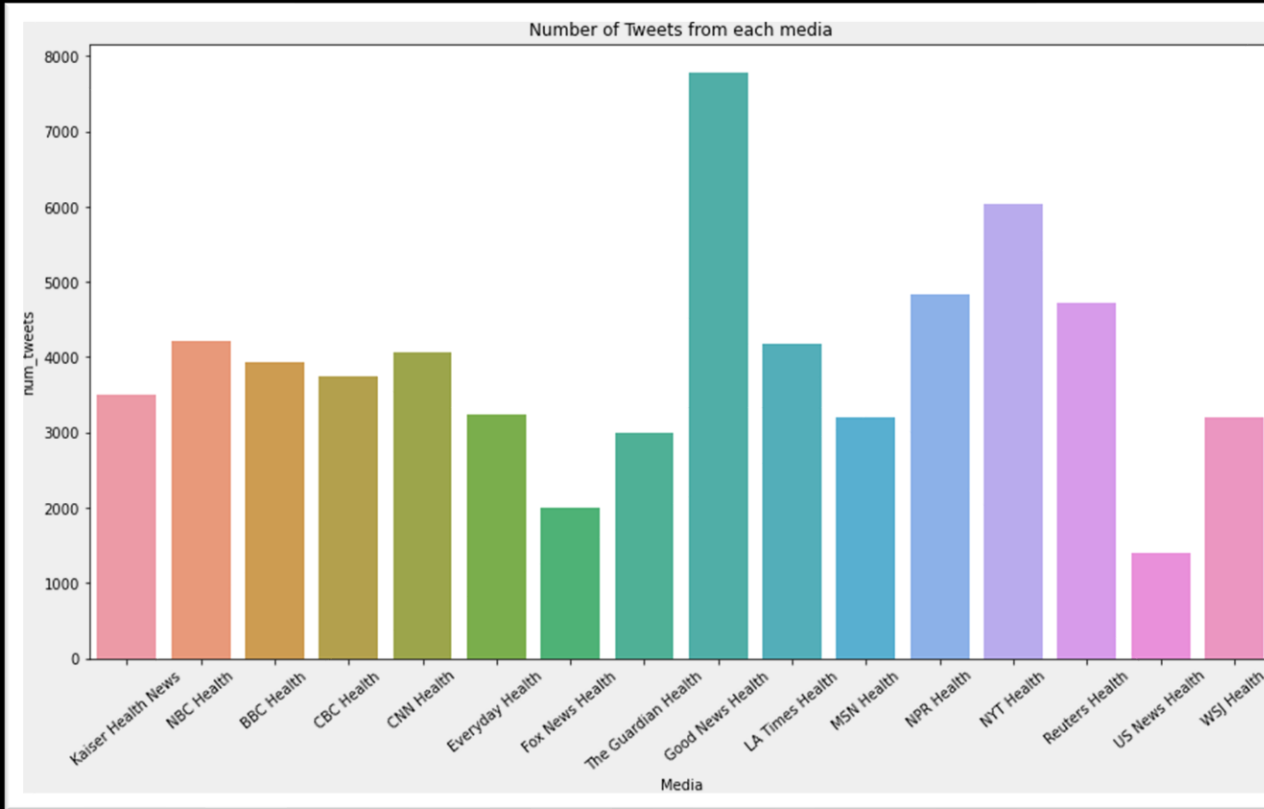
# Health News Tweets Categorization

# Data and Problem

o We saw word vector based model worked well with news articles in module 4

o What if each document is much smaller, specifically Twitter text?

o Can we predict the attributions of tweets or the accounts from texts alone?

o Can we find meaningful latent factors from texts?

o The data was downloaded from UCI Machine Learning Repository

o The data contains 63K tweets from 16 health news accounts

o I will use 1,200 tweets from each media, except for the TFIDF analysis

# Data Distribution

Number of tweets and average words per tweet

# TF/IDF 1-2 ngrams

10 words with the highest TFIDF scores from each news source

Many of them confirm our perceptions

This could explain/predict news sources

Kaiser Health News:
['health law' 'headlines' 'cliff' 'law s' 'kaiserfamfound' 'marketplaces' 'reports today' 'cms' 'reports rt' 'fiscal cliff']
NBC Health:
['safety rules' 'genuine' 'nfl concussion' 'powerful painkiller' 'new safety' 'communion' 'fda proposes' 'bra' 'foster farms' 'november']
BBC Health:
['mps' 'gp' 'ebola video' 'video ebola' 'audio' 'uk ebola' 'care home' 'labour' 'ebola' 'centre']
CBC Health:
['canadian' 'canada s' 'b c' 'canadians' 'p e' 'quebec' 'ebola' 'ford' 'centre' 'facts ebola']
CNN Health:
['rt cnnhealth' 'cnn' 'lost pounds' 'triathlon' 'cnnhealth' 'weightloss' 'lost lbs' 'timehealthland' 'transformation' 'pls']
Everyday Health:
['everydayhealth' 'menshealth' 'jillianmichaels' 'digest' 'rt cspi' 'cspi' 'eatsmartbd' 'psoriasis' 'worldcancerday' 'rt everydayhealth']
Fox News Health:
['disneyland' 'weightloss' 'wakes' 'heart transplant' 'ebola' 'leone' 'sierra leone' 'georgia' 's' 'lodged']
The Guardian Health:
['midwife' 't miss' 'gp' 'day life' 'case missed' 'burnout' 'newsletter' 'andy' 'frontline' 'network s']
Good News Health:
['abs' 'rd' 'weightloss' 'recipes' 'healthyliving' 'fitfluential' 'vacuum' 'greatist' 'gorgeous' 'toned']
LA Times Health:
['cdcgov' 'healthful' 'glutenfree' 'planned parenthood' 'komen' 'rt' 'medi' 'ameracadpeds' 'nejm' 's']
MSN Health:
['clot risk' 'study u' 'gene mutations' 'study' 'heart trouble'  'patients study' 'study study' 's' 'women study' 'blood thinner']
NPR Health:
['npr' 'nprhealth' 'health law' 'ebola' 'health exchanges' 'exchanges' 's' 'insurance website' 'health exchange' 'health']
NYT Health:
['aug' 'briefing' 'op ed' 'nytimes' 'like doctor' 'contributor' 'global health' 'nyt' 'new health' 'letters']
Reuters Health:
['study u' 'mid stage' 'roche' 'eu' 'ebola' 'sanofi' 'abbvie'  'stage study' 'gilead' 'amgen']
US News Health:
['usnews' 'leonardkl' 'rankings' 'rt leonardkl' 'kerigans' 'd love''healthyliving' 'fitness tracker' 'goredforwomen' 'hearthealth']
WSJ Health:
['pharma' 'valeant' 'health law' 'h rt' 'wsj' 'allergan' 'abbvie'  'good morning' 'headlines' 'htt rt']

# TF/IDF Based Clustering And Categorization

- Predictive model using the TF/IDF scores was not successful with 30.18% accuracy score even with the training data
- New tweets can be categorized into "closest" source media
  - It is interesting that the model tells what each tweet "looks like"
  - But it does not accurately predict the tweets from the accounts that were in the training set

- TF/IDF vectors categorizes tweets and source accounts, but not suitable for building predictive models

# TF/IDF and MF for Sources



2 factors and media

The TFIDF vectors are decomposed with 2 components Matrix Factorization

See if the two factors are latent factors identifying the source medias

Looks to me that
- traditional vs unconventional horizontally (f-0)
- Sensational vs dried vertically (f-1)

# Word Count Vector and MF for Tweets

Created word count vectors and decomposed with 2 components Matrix Factorization on each tweet

Does not look like it can distinguish the source accounts, but it does cluster tweets

Decided to build a model to cluster tweets based on word count vector and MF

# Tweet Categorization

- Built a model to categorize tweets into 4 categories
- Used Word Count Vectorization, 5 components Matrix Factorization and Kmeans (k=4)
- Tweets from each source account are categorized as shown in the bar chart

- It gives certain evaluations to the source accounts, and categories to the tweets

- Tweets categorization intuitively make sense

# Evaluation and Conclusion

- Tweet texts are too short to make predictive models
- Tweet analysis should include linked pages, images, and hashtags
- Vectorization/Factorization clustering still yields interesting insights
  - It could give certain labels to twitter accounts
  - It could categorize tweets – could make a basis for recommendation or curation engines
- For validation of unsupervised model, manual review by multiple individuals is mandatory

- I gave my personal interpretations on the 4 categories
  - Category 0: social matters and policies, Category 1: diseases and illness
  - Category 2: personal health such as diet, Category 3: scientific and medical
  - This was not too bad – see what the model said on some tweets today in next 4 slides

**NPR Health News** ✔ @NPRHealth · 12m

Union President John Courtney discusses masking on public transportation

npr.org
Union President John Courtney discusses masking...
NPR's Scott Simon speaks to John Courtney, president of a transit workers' union in California, ...

🔁 1

**NPR Health News** ✔ @NPRHealth · 16h
How Kentucky Republicans blocked all abortions for more than a week

PROTECT SAFE, LEGAL ABORTION
PROTECT SAFE, LEGAL ABORTION

npr.org
How Kentucky Republicans blocked all abortions for more than a week
Even without a Supreme Court ruling, a new Kentucky law shut down abortions for several days before a federal court stepped in. Abortion ...

🔁 1  ♡ 2

## Category 0: Social and Policies

**TIME Health** ✔ @TIMEHealth · Apr 21
Should You Still Wear a Mask on Planes, Trains, and Buses? Here's What the Science Says

time.com
Should You Still Wear a Mask on Planes, Trains and Buses? Here's Wh...
The rules have changed; the virus hasn't. Mask up.

🔁 4  ♡ 2

**TIME Health** ✔ @TIMEHealth · Apr 21
I Tested Positive for COVID-19. Should I Take an Antiviral or Antibody Treatment?

time.com
I Tested Positive for COVID-19. Should I Take a Drug Therapy?
Several medications are available to treat COVID-19, but they're only meant for a specific group of people, and during a very short window ...

♡ 1

**TIME Health** ✓ @TIMEHealth · 22h

At-Home COVID-19 Tests Expire. Here's What to Know About Yours

time.com
At-Home COVID-19 Tests Expire. What to Know About Yours
Here's how long they really last

💬    ⟲    ♡ 1    ⬆️

**Reuters Health** ✓ @Reuters_Health · 15h

Here's what you need to know about the pandemic right now:

reuters.com
What you need to know about the coronavirus right now
Here's what you need to know about the pandemic right now:

💬 1    ⟲    ♡    ⬆️

**WSJ Health** ✓ @WSJhealth · Apr 21

Justice Department unveils charges that range from overcharging for medical services to selling fake vaccination cards

wsj.com
Alleged Covid-19 Fraud Schemes Totaling $150 Million Draw Criminal …
The Justice Department unveiled charges that range from overcharging for medical services to selling fake vaccination cards.

💬    ⟲    ♡ 1    ⬆️

## Category 1: Diseases and Illness

**Everyday Health** ✓ @EverydayHealth · 6h

#Bipolar disorder is widely misunderstood, and it's time to clear up those false beliefs.

everydayhealth.com
7 Myths and Facts About Bipolar Disorder
One sobering fact about bipolar disorder is that as many as half of the people with the condition …

💬    ⟲ 1    ♡ 1    ⬆️

**NPR Health News** ✓ @NPRHealth · 3h

A new puberty guide for kids aims to replace anxiety with self-confidence

npr.org
A new puberty guide for kids aims to replace anxiety with self-confide…
Talking about testes and menses can be super awkward for any kid. A new book tries to take the embarrassment out of growing up – and be …

💬    ⟲ 5    ♡ 7    ⬆️

**WSJ Health** ✓ @WSJhealth · Apr 20

To mask or not to mask: Whether and when to don a mask has largely become a matter of personal choice.

wsj.com
Should You Still Wear a Mask as Federal Mask Mandates Are Dropped…
Whether and when to don a mask has largely become a matter of personal choice.

💬 2    ⟲ 3    ♡ 1    ⬆️

Category 2: Personal Health

# Category 3: Medical Science

**U.S. News Health** ✓ @USNewsHealth · 1h · · ·
These are the Best Children's Hospitals in the nation for neurology and neurosurgery. #BestHospitals



health.usnews.com
These Are the Top Children's Hospitals in the National for Neurology a...
For treating everything from epilepsy to stroke in kids, these medical centers are the best of the best.

💬          ⟲ 1          ♡ 2          ⬆

**WSJ Health** ✓ @WSJhealth · Apr 21 · · ·
Amid rising marijuana use in the U.S., researchers are exploring risks to bystanders, children—and pets



wsj.com
Rising Marijuana Use Presents Risks to Pets, Bystanders
A variety of recent studies have examined the incidental effects of pot use, including one tying its legalization to increased cannabis ...

💬          ⟲ 2          ♡ 5          ⬆

**BBC Health News** ✓ @bbchealth · 20h · · ·
Covid cases in Scotland fall by 30,000 in a week



bbc.com
Covid cases in Scotland fall by 30,000 in a week
It is the fourth week in a row the number of positives tests has dropped, according to official figures.

💬 1          ⟲ 1          ♡ 6          ⬆

**BBC Health News** ✓ @bbchealth · 19h · · ·
Covid-19: Three more deaths and 390 in hospital



bbc.com
Covid-19: Three more deaths and 390 in hospital
The total number of deaths linked to Covid-19 in NI since the start of the pandemic is 3,405.

💬 1          ⟲ 1          ♡          ⬆

**U.S. News Health** ✓ @USNewsHealth · 2h · · ·
.@TexasChildrens is the top-ranked hospital in treating high-risk heart defects and helping children with severe congenital heart disease.



health.usnews.com
Best Hospitals for Children With Severe Congenital Heart Disease
These 26 hospitals are top-ranked in pediatric heart care and experienced in treating high-risk defects.

💬          ⟲          ♡ 1          ⬆