



DSE 203 - Knowledge Graph Report

NFL Players with Chronic Traumatic Encephalopathy (CTE)

Table of Contents

Section 1: Overview	3
1.1 Context.....	3
1.2 Objective	3
1.3 Scope of Work	3
Section 2: Data Sources.....	4
Section 3: Data Sources 1 & 2 Integration Strategy	6
Section 4: Data Source 3 Named Entity Recognition (NER) & LDA Topic Modeling	13
Section 5: Data Sources 2 & 3 Integration Strategy	17
Section 6 Knowledge Graph Implementation Strategy	19
Section 7 Analytical Queries	19
Section 8 Project Files	20
Section 9 Revision History.....	20

Section 1: Overview

1.1 Context

Many former National Football League (NFL) players have been diagnosed with or have had chronic traumatic encephalopathy, or CTE. A definitive diagnosis so far can be made only post-mortem. However, an increasing number of former players are reporting symptoms of CTE.

According to 2017 study on brains of deceased gridiron football players, 99% of tested brains of NFL players, 88% of Canadian Football League (CFL) players, 64% of semi-professional players, 91% of college football players, and 21% of high school football players had various stages of CTE. However, this study had several limitations, including possible selection bias as families of players with symptoms of CTE are far more likely to donate brains to research than those without signs of the disease. Despite the limitations, the study still showed that CTE is far more common than once believed.

1.2 Objective

Organize data from multiple sources, about entities of interest associated with CTE in the NFL, and forge connections between them to represent a network of real-world entities—i.e. objects, events, situations, or concepts—and illustrates the relationship between them in a graph database and visualized as a graph model.

1.3 Scope of Work

Integrate data from three different sources which combines a mixture of structured, semi-structured, and unstructured data to construct a knowledge graph per below provided project requirements for UCSD's DSE203 final project.

Project Requirements: *Your knowledge graph need not be humongous, but it should have a reasonable size (nodes + edges)*

- *100 is too low, 100k is high*
- *The number of edge labels should not be less than 7*
- *The number of node labels should not be less than 7*

Your methods should include some degree of information extraction from text and some degree of value matching

Try to make your use cases realistic. Your example queries on the knowledge graph

- *Should have some degree of complexity*
- *May include some analytical operations on graphs that are permitted by Neo4J (and Python if needed)*

Section 2: Data Sources

The project team combined a mixture of datasets from the following publicly available data sources.

Data source 1: Kaggle - NFL Statistics

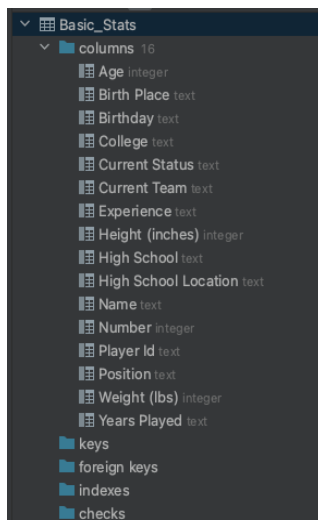
****Note:** The project team only utilized the Basic_Stats.csv dataset from data-source-1.

Data Variety - Structured data.

Data Source URL: <https://www.kaggle.com/datasets/kendallgillies/nflstatistics>

Download URL: <https://github.com/mona-jandro-camm/dse203/tree/main/Datasets>

DDL Schema:



The image shows a DDL Schema for a table named 'Basic_Stats'. It lists 16 columns with their data types: Age (integer), Birth Place (text), Birthday (text), College (text), Current Status (text), Current Team (text), Experience (text), Height (inches) (integer), High School (text), High School Location (text), Name (text), Number (integer), Player Id (text), Position (text), Weight (lbs) (integer), and Years Played (text). Below the columns, there are sections for keys, foreign keys, indexes, and checks, all of which are currently empty.

	"Player Id"	"Name"	"Age"	"Birthday"	"College"	"Current Status"	"Current Team"
1	bobbyfowler/2514295	Fowler, Bobby	56	9/11/1960	Louisiana Tech	Retired	<null>
2	quinnjohnson/79593	Johnson, Quinn	30	9/30/1986	LSU	Retired	<null>
3	l.t.walton/2552444	Walton, L.T.	25	3/31/1992	Central Michigan	Active	Pittsburgh Steelers
4	jordanleggett/2557885	Leggett, Jordan	22	1/31/1995	Clemson	Active	New York Jets
5	georgegonda/2515030	Gonda, George	98	2/23/1919	Duquesne	Retired	<null>
6	floydraglin/2523695	Raglin, Floyd	56	2/10/1961	Southern University	Retired	<null>
7	jaredzabransky/2495791	Zabransky, Jared	33	12/4/1983	Boise State	Retired	<null>

Data source 2: Wikipedia - NFL players with chronic traumatic encephalopathy

Data Variety - Semi structured and unstructured data.

Data Source URL:

https://en.wikipedia.org/wiki/List_of_NFL_players_with_chronic_traumatic_encephalopathy



Main page
Contents
Current events
Random article
About Wikipedia
Contact us
Donate

Contribute
Help
Learn to edit
Community portal
Recent changes
Upload file
Tools
What links here
Related changes
Special pages
Permanent link

Article Talk

List of NFL players with chronic traumatic encephalopathy

From Wikipedia, the free encyclopedia

A large number of former **National Football League** (NFL) players have been diagnosed with or have had **chronic traumatic encephalopathy**, or CTE. are reporting symptoms of CTE.

According to 2017 study on brains of deceased **gridiron football** players, 99% of tested brains of NFL players, 88% of **Canadian Football League** (CF players had various stages of CTE. However, this study had several limitations, including possible selection bias as families of players with symptom limitations, the study still showed that CTE is far more common than once believed.^[1]

Other common injuries include injuries of legs, arms, and lower back.^{[2][3][4][5]}

Contents [hide]

- 1 Players affected
- 2 Former players with CTE confirmed post-mortem
- 3 Deceased players suspected of having had CTE
- 4 Living former players diagnosed with CTE or ALS or reporting symptoms consistent with CTE or ALS
- 5 Former players listed as plaintiffs in lawsuits against the NFL for concussion-related injuries received after playing
- 6 See also
- 7 References

Former players with CTE confirmed post-mortem [edit]

A common definitive test currently can be made only by examining the brain tissue of a deceased victim.

As the families of many deceased players wish to keep their medical information private, the following list is incomplete. A brain injury study conducted at the **Boston** signs of CTE,^[13] and additional players have so far been confirmed with CTE separately. A new list released in November 2016 mentions CTE in 90 of 94 brains of examined showed signs of CTE.^[15]

- | | | | | |
|--|---|--|---|---|
| • Phillip Adams ^[16] | • Dwight Clark ^[25] | • Ray Easterling ^{[31][8][32]} | • Wally Hilgenberg ^[40] | • Jim Kiick ^{[21][22]} |
| • George Andrie ^[17] | • Greg Clark ^[26] | • Grant Feasel ^[33] | • Glen Ray Hines ^[41] | • Terry Long ^[13] |
| • Jovan Belcher ^[18] | • Daniel Colchico ^[27] | • Frank Gifford ^{[31][34]} | • Jim Hudson ^[42] | • Rob Lytle ^[45] |
| • Forrest Blue ^[19] | • Lou Creekmur ^[28] | • Cookie Gilchrist ^{[31][35]} | • Gerry Huth ^[27] | • John Mackey ^{[31][46]} |
| • Colt Brennan ^[20] | • Art DeCarlo ^[27] | • John Gimsley ^[36] | • Vincent Jackson ^[43] | • Ollie Matson ^{[31][47]} |
| • Nick Buoniconti ^{[21][22]} | • Shane Dronett ^[29] | • Jason Hairston ^[37] | • John Henry Johnson ^[44] | • Tom McHale ^{[31][48]} |
| • Lew Carpenter ^[23] | • Dave Duerson ^{[30][31]} | • Chris Henry ^{[31][38]} | • Tom Keating ^[27] | • Earl Morrall ^{[22][49]} |
| • Ronnie Caveness ^[24] | • Pete Duranko ^[27] | • Aaron Hernandez ^[39] | • Bob Kuechenberg ^[22] | • Larry Morris ^[50] |

Data source 3: UCSD's AWESOME Database - US Newspaper Articles

Data Variety- Structured and unstructured data.

Table name: usnewspaper

Connection details:

host - awesome-hw.sdsc.edu

database name - postgres

DDL Schema:

usnewspaper	
columns 11	
news	text
id	integer = nextval('usnewspap...)
collectiondate	date
title	varchar(600)
url	varchar(600)
publishdate	date
author	text[]
keywords	text[]
src	varchar(400)
language	varchar(2)
newsindex	tsvector
keys 1	
usnewspaper_pkey	(id)
foreign keys	
indexes 3	
usnewspaper_pkey	(id) UNIQUE
usnewspaper_publishdate_index	(publishdate)
usnewspaper_src_index	(src)

	title	news	keywords
1	Q&A: Here's what jury must con...	Following a month of wildly dispar...	{matt,jury,lawsuit,testified,gee,us
2	Former NFL WR Demaryius Thomas...	Demaryius Thomas had stage 2 CTE a...	{wr,struggle,cte,family,stage,trauma
3	Brian Urlacher says some ex-NF...	Hall of Fame linebacker Brian Url...	{cte,lawsuit,theres,brian,players,f
4	Former NFL receiver, 33, had s...	Demaryius Thomas had stage 2 CTE a...	{encephalopathy,mckee,dont,death,sy
5	Ashley Massaro wanted to donat...	(CNN) After former WWE star Ashley...	{wrestler,donate,concussion,wrestli
6	TorHoerman Law Files Lawsuit A...	CountryUnited States of America ...	{islands,law,states,kingdom,peoples
7	ASU, NCAA sued over death of L...	The father of a former Arizona Sta...	{ncaa,suffered,disease,franklin,dea

Section 3: Data Sources 1 & 2 Integration Strategy

Method:

Extract NFL CTE affected player names from Wikipedia page and perform player name matching with NFL basic statistics dataset obtained from Kaggle. To extract the Wikipedia player names, the project team utilized Python's BeautifulSoup library, and for the name matching between both datasets, py_entitymatching's attribute blocker EM process was utilized.

Technology Stack:

The following technology stack was primarily used to integrate data sources 1 & 2.

- Python (Jupyter Notebooks)
- SQL (Postgres SQL, for Analysis only)
- Python Libraries

```

import py_stringmatching as sm
import py_entitymatching as em
import pandas as pd
import numpy as np
import re, string, math, time
import wikipedia
import stanza
import requests
import csv
from bs4 import BeautifulSoup
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

```

Data Integration Steps:

1. Load structured CSV dataset from data source #1 - Kaggle into a DataFrame via `py_entitymatching's read_csv_metadata` method.

Sample Code Snippet:

Load Kaggle Dataset

```
basic_stats_df = em.read_csv_metadata("./Datasets/Basic_Stats.csv", key='Player Id')
```

2. Extract semi-structured and unstructured data from data source #2 - Wikipedia. The data extraction was done via Python's Wikipedia and BeautifulSoup libraries.
***Note: BeautifulSoup provides additional drill-down functionality to get at specific HTML tags in comparison to Python's Wikipedia library.*

Sample Code Snippet:

Extract Wikipedia Data

```

: wiki_title = 'List of NFL players with chronic traumatic encephalopathy'
  wiki_url = 'https://en.wikipedia.org/wiki/List_of_NFL_players_with_chronic_traumatic_encephalopathy'

# Python Wikipedia library
wiki_page_object = wikipedia.page(wiki_title)

# Python BeautifulSoup
wiki_page = requests.get(wiki_url)
soup = BeautifulSoup(wiki_page.content, "lxml")

```

3. Perform Named Entity Recognition (NER) and extract "PERSON" entities with the help of Python's Stanza library for data source #2 - Wikipedia. The "PERSON" entity data was placed into separate data/lists structures based on each header category they are listed in data source #2. Those five different "PERSON" header categories are listed below.
 - Players Affected Wiki Section

- *Former Players affected with CTE Wiki Section*
- *Deceased players suspected of having had CTE Wiki Section*
- *Living former players diagnosed with CTE or ALS or reporting symptoms consistent with CTE or ALS Wiki Section*
- *Former players listed as plaintiffs in lawsuits against the NFL for concussion-related injuries received after Wiki playing Section*

Sample Code Snippet:

Stanza - stanford NLP

```
nlp = stanza.Pipeline('en', processors='tokenize,mwt,ner', use_gpu=False, pos_batch_size=3000, download_method=None)
```

Process Former Players affected with CTE Wiki Section

```
# Lists to store player names by category
former_players_post_mortem_ls = []
pm_former_players_ls = []

# Wiki-Extract Former players with CTE confirmed post-mortem
results = soup.select('ul')[1]
former_players_post_mortem_ls = results.find_all("a")

# Set start time to calculate compute time
start_time = time.time()

# PASS-1: Compute NER with wiki page links - Former players with CTE confirmed post-mortem
for index in former_players_post_mortem_ls:
    doc = nlp(str(index))

    # Extract PERSON entities
    for ent in doc.ents:
        if (ent.type == 'PERSON'):
            clean_name = re.split('</a', ent.text)[0]
            pm_former_players_ls.append(clean_name)

# Dedupe list contents
pm_former_players_ls = [*set(pm_former_players_ls)]

print("Exec time --- %s seconds ---" % (time.time() - start_time))
print(f'# of person: {len(pm_former_players_ls)}')
```

Exec time --- 19.32221007347107 seconds ---
of person: 63

4. Text normalization and pre-processing to clean player names, lower and remove punctuations from text.

Sample Code Snippet:

Text Normalization & Preprocessing

```
# -----
# Normalize player "Name" in Kaggle basic stats
# -----
basic_stats_df['Clean_Name'] = basic_stats_df.Name.str.lower().map(lambda s: s.split()[1] + ' ' + s.split()[0]).replace('[^\w\s]',' ', regex=T

# -----
# Remove punctuations & lower name
# -----
def remove_punc(name):
    punc = '!()-[]{};:'"\, <>./?@$%^&*~'"
    for ele in name:
        if ele in punc:
            name = name.replace(ele, " ")
    return name.lower().strip()

affected_players_ls = [remove_punc(i) for i in affected_players_ls]
pm_former_players_ls = [remove_punc(i) for i in pm_former_players_ls]
suspected_deceased_players_ls = [remove_punc(i) for i in suspected_deceased_players_ls]
cte_als_former_players_ls = [remove_punc(i) for i in cte_als_former_players_ls]
players_nfl_lawsuits_ls = [remove_punc(i) for i in players_nfl_lawsuits_ls]
```

- Combine the five different header categories data structures which contain NFL player names from data source #2 into one single DataFrame while creating following groups for each header category.

Wiki Header Section	CTE Category Group
Players Affected Wiki Section	affected_players
Former Players affected with CTE Wiki Section	pm_former_players
Deceased players suspected of having had CTE Wiki Sction	suspected_deceased_players
Living former players diagnosed with CTE or ALS or reporting symptoms consistent with CTE or ALS Wiki Section	cte_als_former_players
Former players listed as plaintiffs in lawsuits against the NFL for concussion-related injuries received after Wiki playing Section	players_nfl_lawsuits

Sample Code Snippet:

```
# Combine dataframes
frames = [affected_players_df, pm_former_players_df, suspected_deceased_players_df, cte_als_former_players_df, players_nfl_lawsuits_df]
wiki_cte_players_df = pd.concat(frames)
wiki_cte_players_df
```

	cte_category	Clean_Name
0	affected_players	larry johnson
1	affected_players	stabler
2	affected_players	busm
3	affected_players	ken stabler
4	affected_players	johnson
...
1875	players_nfl_lawsuits	richard cash
1876	players_nfl_lawsuits	ricky siglar
1877	players_nfl_lawsuits	john turner
1878	players_nfl_lawsuits	roderick coleman
1879	players_nfl_lawsuits	david alexander

- Perform entity matching Attribute Blocking on data source #1 & #2 player names and create a new DataFrame.

Sample Code Snippet:

Block DataFrames to get Candidate set

```
# Instantiate blocker objects:
# -----
# Create attribute equivalence blocker
ab = em.AttrEquivalenceBlocker()

# Create overlap blocker
ob = em.OverlapBlocker()
```

ii. Attribute Block by 'player_name'

***** BETTER RESULTS THAN OVERLAP BLOCK *****

```
# Block using 'full_name_dob' attribute
ab_fullname_cand = ab.block_tables(basic_stats_df, wiki_person_df, 'Clean_Name', 'Clean_Name', allow_missing=False,
                                   l_output_attrs=['Player Id', 'Name', 'Age', 'Current Status', 'Birthday', 'College', 'High School', 'Clean_Nam
                                   r_output_attrs=['rec_id', 'Clean_Name', 'cte_category'], n_jobs=2)
```

```
# Distinct matched candidates - Kaggle vs. Wiki page
ab_fullname_cand.groupby("ltable_Player Id").first().to_csv('./AB_names_matched.csv')
ab_fullname_cand.groupby("ltable_Player Id").first()
```

	_id	rtable_rec_id	ltable_Name	ltable_Age	ltable_Current Status	ltable_Birthday	ltable_College	ltable_High School	ltable_Clean_Name	rtable_Clean
ltable_Player Id										
aaronbeasley/2499587	697	1731	Beasley, Aaron	43.0	Retired	7/7/1973	West Virginia	None	aaron beasley	aaron
aaronjones/2558116	16	285	Jones, Aaron	22.0	Active	12/2/1994	Texas-El Paso	Burges HS	aaron jones	aarc
adamhaayer/2504632	705	1353	Haayer, Adam	40.0	Retired	2/22/1977	Minnesota	None	adam haayer	adarr
adriandingle/2500398	491	1528	Dingle, Adrian	39.0	Retired	6/25/1977	Clemson	None	adrian dingle	adria

- Perform Named Entity Recognition (NER) and extract "ORG" entities with the help of Python's Stanza library for data source #2 - Wikipedia. The "ORG" entity data was placed into separate DataFrame in order to map NFL player names to an organization entity.

Sample Code Snippet:

```
# PASS-1: Compute NER with wiki page to extract ORGs
doc = nlp(NormalizeText(wiki_page_object.content))

# Extract ORG entities
for ent in doc.ents:
    if (ent.type == 'ORG'):
        wiki_org_ls.append(ent.text)

# Dedupe list contents
wiki_org_ls = [*set(wiki_org_ls)]

print("Exec time --- %s seconds ---" % (time.time() - start_time))
print(f'# of organizations: {len(wiki_org_ls)}')
```

```
Exec time --- 5.689316749572754 seconds ---
# of organizations: 14
```

```
# Display Wiki page Organizations
wiki_org_df = pd.DataFrame(data= wiki_org_ls, columns=['wiki_org_name'])
wiki_org_df['ord_id'] = wiki_org_df.index+1
wiki_org_df
```

	wiki_org_name	ord_id
0	Atlanta Falcons	1
1	Canadian Football League CFL	2
2	CTE The Brain Bank	3
3	Colts	4
4	The Boston University School Medicine BUSM	5
5	National Institute Occupational Safety Health NIOSH	6
6	NFL All Pro	7
7	NFL	8

8. The blocked/matched NFL player names from data sources #1 & #2 are mapped to their parent “ORG” entity, i.e. NFL, and the graph model direction, i.e. parent-to-child, attribute is added to DataFrame created in step #6.

Sample Code Snippet:

Assign Parent Node & Direction to Players DataFrame

```
ab_fullname_cand['parent'] = [wiki_org_df.query("wiki_org_name == 'NFL'")['wiki_org_name'].values[0]] * len(ab_fullname_cand)
ab_fullname_cand['direction'] = ['parent_to_child'] * len(ab_fullname_cand)
ab_fullname_cand['ltable_Age'] = pd.to_numeric(ab_fullname_cand.ltable_Age, downcast='integer')
ab_fullname_cand
```

Itable_Name	Itable_Age	Itable_Current Status	Itable_Birthday	Itable_College	Itable_High School	Itable_Clean_Name	rtable_Clean_Name	rtable_cte_category	parent	direction
Verdin, Clarence	53.0	Retired	6/14/1963	Louisiana-Lafayette	NaN	clarence verdin	clarence verdin	players_nfl_lawsuits	NFL	parent_to_child
Lewis, Kevin	37.0	Retired	4/26/1980	Virginia Tech	NaN	kevin lewis	kevin lewis	players_nfl_lawsuits	NFL	parent_to_child
Lewis, Kevin	50.0	Retired	11/14/1966	Northwestern State-Louisiana	NaN	kevin lewis	kevin lewis	players_nfl_lawsuits	NFL	parent_to_child
Lewis, Kevin	38.0	Retired	10/6/1978	Duke	NaN	kevin lewis	kevin lewis	players_nfl_lawsuits	NFL	parent_to_child
Chambers, Chris	38.0	Retired	8/12/1978	Wisconsin	NaN	chris chambers	chris chambers	players_nfl_lawsuits	NFL	parent_to_child
...
Roberts, Walter	75.0	Retired	2/15/1942	San Jose State	NaN	walter roberts	walter roberts	players_nfl_lawsuits	NFL	parent_to_child
Anderson, Donny	74.0	Retired	5/16/1943	Texas Tech	NaN	donny anderson	donny anderson	players_nfl_lawsuits	NFL	parent_to_child
Graff, Neil	67.0	Retired	1/12/1950	Wisconsin	NaN	neil graff	neil graff	players_nfl_lawsuits	NFL	parent_to_child
Campbell, Lamar	40.0	Retired	8/29/1976	Wisconsin	NaN	lamar campbell	lamar campbell	players_nfl_lawsuits	NFL	parent_to_child
Welter, Tom	53.0	Retired	2/24/1964	Nebraska	NaN	tom welter	tom welter	players_nfl_lawsuits	NFL	parent_to_child

9. Construct and unit-test Neo4j knowledge graph to for data source #1 & #2.

Basic Rules logic: The nodes are simply players and their parent “ORG” entity NFL.

The edges are the five different CTE header categories from data source #2.

Sample Code Snippet:

Construct Neo4j Node CSV File

```
def processNodes(data, node_file):
    nodes = {}
    counter = 1
    node_header = [":ID", "Name", "PlayerID", "Age", "Birthday", "Status", "College", ":LABEL"]

    # Set start time to calculate compute time
    start_time = time.time()

    # Construct node map:
    for index, row in data.iterrows():
        parent_node_id = row.parent
        child_node_id = row['table_Player ID']

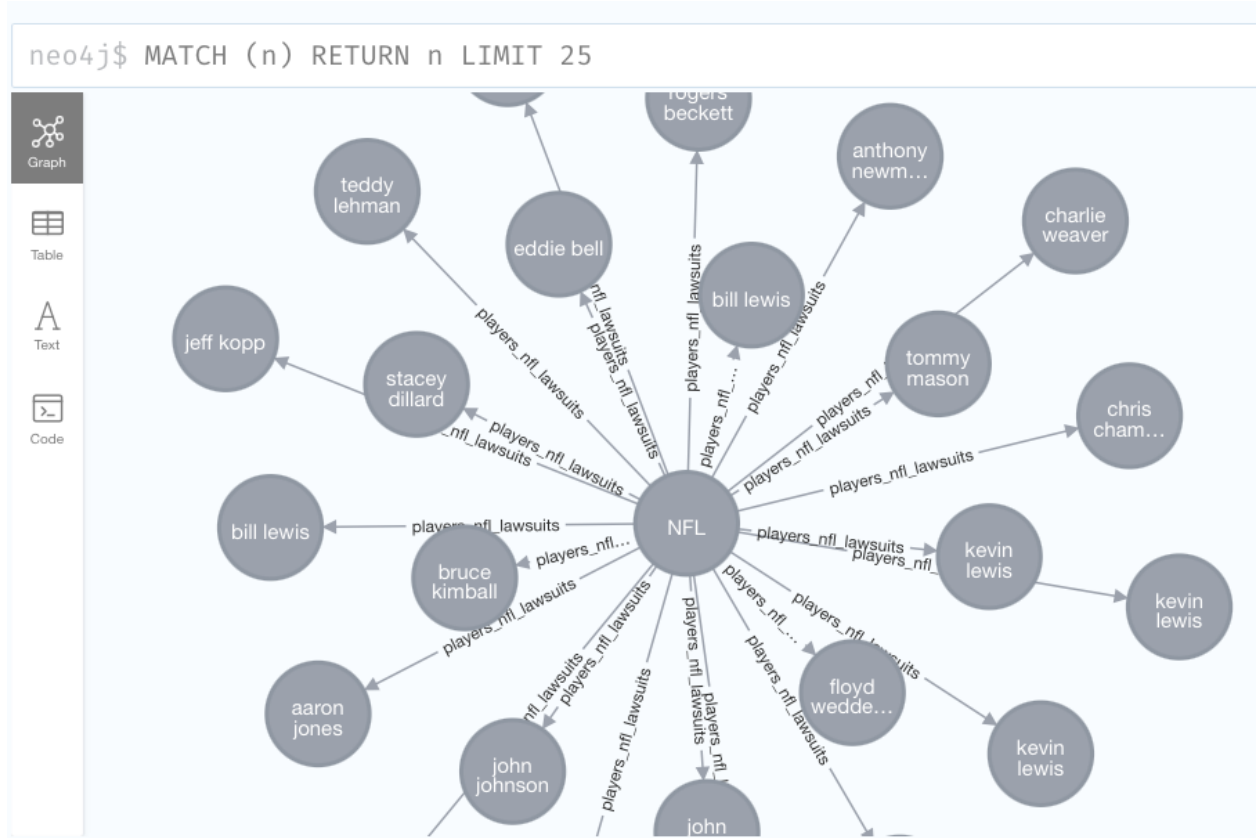
        if parent_node_id is None or child_node_id is None:
            continue;

        # Check if parent node already mapped, otherwise add
        if not bool([i for i in nodes if nodes[i][0] == parent_node_id]):
            nodes[counter] = [parent_node_id, parent_node_id, '', '', '', parent_node_id]
            counter+=1

        # Check if child node already mapped, otherwise add
        if not bool([i for i in nodes if nodes[i][0] == child_node_id]):
            nodes[counter] = [row['table_Clean_Name'] if child_node_id != 'NFL' else ''
                             , child_node_id if child_node_id != 'NFL' else ''
                             , row['table_Age'] if (child_node_id != 'NFL' and row['table_Age'] == row['table_Age']) else ''
                             , row['table_Birthday'] if (child_node_id != 'NFL' and row['table_Birthday'] == row['table_Birthday']) else ''
                             , row['table_Current Status'] if child_node_id != 'NFL' else ''
                             , row['table_College'] if child_node_id != 'NFL' else ''
                             , row['table_Clean_Name'] if child_node_id != 'NFL' else ''
            counter+=1

    # write nodes CSV file
    with open(node_file, 'w', newline='') as f:
        writer = csv.writer(f)
        writer.writerow(node_header)
        for node in nodes:
            if nodes[node][0] == 'NFL':
                writer.writerow([node, nodes[node][0], nodes[node][0], nodes[node][0], nodes[node][0], nodes[node][0], nodes[node][0], nodes[node][0]])
            else:
                writer.writerow([node, nodes[node][0], nodes[node][1], nodes[node][2], nodes[node][3], nodes[node][4], nodes[node][5], nodes[node][6]])

    # compute execution time
    exec_time = time.time() - start_time
```



Section 4: Data Source 3 Named Entity Recognition (NER) & LDA Topic Modeling

Method:

Extract NFL CTE affected player names and keywords from USNewspaper articles using named entity recognition (NER). The project team created a LDA topic model in order to build further knowledge about each player from the unstructured data available in the newspaper articles. The team utilized various NLP Python libraries as outlined below including Spacy for NER and gensim for the LDA topic modeling.

Technology Stack:

The following technology stack was primarily used to extract information from data source 3.

- Python (Jupyter Notebooks)
- SQL (Postgres SQL)
- Python Libraries

```

import sqlalchemy as sal
from sqlalchemy import text

import nltk

import gensim
import gensim.corpora as corpora
from gensim.utils import simple_preprocess
from gensim.models import CoherenceModel

import pyLDAvis
import pyLDAvis.gensim_models as gensimvis
import matplotlib.pyplot as plt

import spacy

```

Steps:

1. Load unstructured dataset from data source #2 - USnewspapers into a DataFrame using sqlalchemy and Pandas.

```

1 engine = sal.create_engine(
2     "postgresql+psycopg2://ag_class:WUcgdfQl@awesome-hw.sdsc.edu/postgres"
3 )
4 conn = engine.connect()

1 sql_query = text(
2     """SELECT DISTINCT title, news, keywords
3         FROM usnewspaper
4         WHERE ARRAY['cte', 'lawsuit']::text[] <@ keywords and news is not null
5     UNION
6     SELECT DISTINCT title, news, keywords
7         FROM usnewspaper
8         WHERE ARRAY['nfl', 'helmet']::text[] <@ keywords and news is not null
9     UNION
10    SELECT DISTINCT title, news, keywords
11        FROM usnewspaper
12        WHERE ARRAY['nfl', 'brain']::text[] <@ keywords and news is not null
13    UNION
14    SELECT DISTINCT title, news, keywords
15        FROM usnewspaper
16        WHERE ARRAY['encephalopathy']::text[] <@ keywords AND news is not null;"""
17 )
18 result = conn.execute(sql_query)

```

```
1 data = [i for i in result]
```

```
1 df = pd.DataFrame(data, columns=["title", "news", "keywords"])
```

2. Use Spacy to extract named entities of type "ORG" and "PERSON".

```
1 nlp_spacy = spacy.load("en_core_web_sm")
```

```
1 docs = list(nlp_spacy.pipe(df["news"]))
```

```

1 list_of_ents = []
2 for doc in docs:
3     list_of_ents.append(
4         list(set([ent.text for ent in doc.ents if (ent.label_ == "ORG") or (ent.label_ == "PERSON")]))
5     )

```

```
1 df["named_entities"] = list_of_ents
```

3. Prepare the dataset for LDA topic model using a combination of re, Spacy, and gensim.

```
1 stop_words = stopwords.words("english")
```

```
1 # Convert to list
2 data = df.news.values.tolist()
3
4 # Remove Emails
5 data = [re.sub(r"\S*\S*\s?", "", sent) for sent in data]
6
7 # Remove new line characters
8 data = [re.sub(r"\s+", " ", sent) for sent in data]
9
10 # Remove distracting single quotes
11 data = [re.sub("'", "", sent) for sent in data]
12 data = [re.sub("`", "", sent) for sent in data]
13 data = [re.sub("-", "", sent) for sent in data]
14
15 #print(data[:1])
```

```
1 def sent_to_words(sentences):
2     for sentence in sentences:
3         yield (
4             gensim.utils.simple_preprocess(str(sentence), deacc=True)
5             ) # deacc=True removes punctuations
6
7
8 data_words = list(sent_to_words(data))
9
10 # print(data_words[:1])
```

```
1 def sent_to_words(sentences):
2     for sentence in sentences:
3         yield (
4             gensim.utils.simple_preprocess(str(sentence), deacc=True)
5             ) # deacc=True removes punctuations
6
7
8 data_words = list(sent_to_words(data))
9
10 # print(data_words[:1])
```

```
1 # Define functions for stopwords and lemmatization
2 def remove_stopwords(texts):
3     return [[word for word in doc if word not in stop_words] for doc in texts]
4
5
6 def lemmatization(texts, allowed_postags=["NOUN", "ADJ", "VERB", "ADV"]):
7     """https://spacy.io/api/annotation"""
8     texts_out = []
9     for sent in texts:
10         doc = nlp(" ".join(sent))
11         texts_out.append(
12             [token.lemma_ for token in doc if token.pos_ in allowed_postags]
13         )
14     return texts_out
15
16
17 def bigrams_and_trigrams(texts):
18
19     # Add bigrams and trigrams to docs (only ones that appear 2 times or more).
20     bigram = gensim.models.Phrases(texts, min_count=2)
21     for idx in range(len(texts)):
22         for token in bigram[texts[idx]]:
23             if " " in token:
24                 # Token is a bigram, add to document.
25                 texts[idx].append(token)
26     return texts
```

```

1 # Remove Stop Words
2 data_words_nostops = remove_stopwords(data_words)
3
4 # Initialize spacy 'en' model, keeping only tagger component (for efficiency)
5 # python3 -m spacy download en
6 nlp = spacy.load("en_core_web_sm", disable=["parser", "ner"])
7
8 # Do lemmatization keeping only noun, adj, vb, adv
9 data_lemmatized = lemmatization(
10     data_words_nostops, allowed_postags=["NOUN", "ADJ", "VERB", "ADV"]
11 )
12
13 data_bigrams = bigrams_and_trigrams(data_lemmatized)
14 data_trigrams = bigrams_and_trigrams(data_bigrams)
15 # print(data_lemmatized[0])
16 # print(data_bigrams[0])
17 # print(data_trigrams[0])

```

```

1 # Create Dictionary
2 id2word = corpora.Dictionary(data_lemmatized)
3
4 # Create Corpus
5 texts = data_lemmatized
6
7 # Term Frequency list
8 corpus = [id2word.doc2bow(text) for text in texts]
9
10 # View
11 # print(corpus[0])

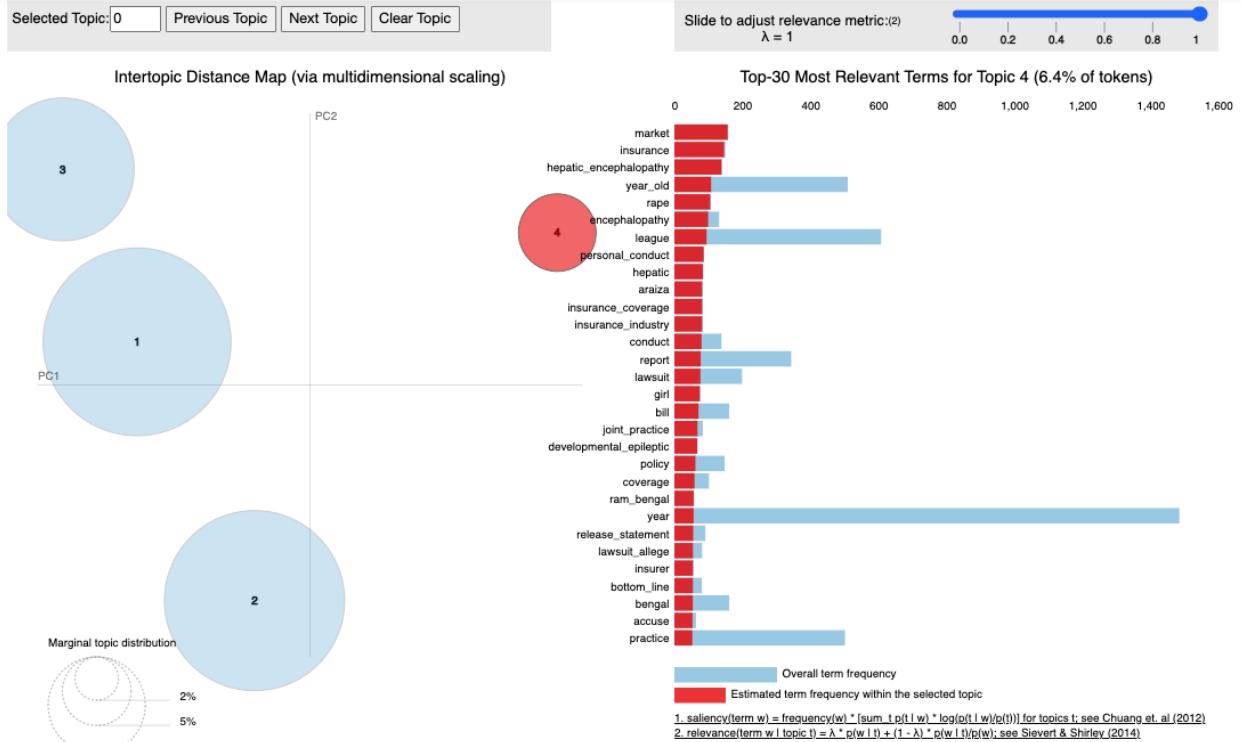
```

4. Build LDA topic model with 4 topics and display it, focusing on one topic of interest in particular, that is, a topic with words associated with “personal conduct” issues.

```

1 # Build LDA model
2 lda_model = gensim.models.ldamodel.LdaModel(
3     corpus=corpus,
4     id2word=id2word,
5     num_topics=4,
6     random_state=100,
7     update_every=1,
8     chunksize=100,
9     passes=10,
10     alpha="auto",
11     per_word_topics=True,
12 )

```



Section 5: Data Sources 2 & 3 Integration Strategy

We found articles, where the list of players from data source 1 and 2 were mentioned and we created an additional set of nodes and relationships for the knowledge graph.

```
[ ]: df_nes_exploded = df.explode("named_entities")

df_names_matched = df_nes_exploded[
    df_nes_exploded["named_entities"]
        .str.lower()
        .isin(ab_fullname_cand["ltable_Clean_Name"])
]
df_names_matched = df_names_matched.loc[
    df_names_matched["title"].drop_duplicates().index
]
df_names_matched
```

```
[ ]: def processNewsRelations(news_data, player_nodes, rel_file):
    relation_header = [":START_ID", ":END_ID", ":TYPE"]
    relation_data = []
    counter = 733

    # Construct relation map:
    for index, row in news_data.iterrows():
        player_id = [
            p_id
            for p_id, val in player_nodes.items()
            if val[0] == row["named_entities"].lower()
        ]
        if len(player_id) == 1:
            relation_data.append([player_id[0], counter, "MENTIONED_IN"])
            counter += 1

    # write relation file
    with open(rel_file, "w", newline="") as f:
        writer = csv.writer(f)
        writer.writerow(relation_header)
        for row in relation_data:
            writer.writerow(row)
```

```
[ ]: def processNewsNodes(data, node_file):
    nodes = {}
    counter = 732
    node_header = ["newsId:ID", "Title", "Publish Date", "Source", "URL", ":LABEL"]

    # Construct node map:
    for index, row in data.iterrows():

        counter += 1
        nodes[counter] = [
            row["title"],
            row["publishdate"],
            row["src"],
            row["url"].strip(),
            f"{row['src'].split('.')[1]}_{row['publishdate']}",
        ]

    # write nodes CSV file
    with open(node_file, "w", newline="") as f:
        writer = csv.writer(f)
        writer.writerow(node_header)
        for node in nodes:
            writer.writerow(
                [
                    node,
                    nodes[node][0],
                    nodes[node][1],
                    nodes[node][2],
                    nodes[node][3],
                ]
            )
```

Section 6 Knowledge Graph Implementation Strategy

Used command line to import node and relationship CSV files that were created in the Python notebook.

Import node command:

```
./bin/neo4j-admin import --force --multiline-fields=true --nodes=./import/CTE_Nodes.csv --  
nodes=./import/CTE_NEWS_Nodes.csv --relationships=./import/CTE_Relations.csv --  
relationships=./import/CTE_NEWS_Relations.csv
```

```
[ ]: # Set data dump path for Neo4j  
# neo4j_data_path = "/Users/camm/Library/NEO4J_HOME/import"  
neo4j_data_path = "/Users/galore/Downloads/neo4j-community-4.4.14/import/"  
  
# Construct Node CSV file  
news_node_map = processNewsNodes(  
    df_names_matched, f"{neo4j_data_path}/CTE_NEWS_Nodes.csv"  
)  
node_map, exec_time = processNodes(  
    ab_fullname_cand.copy(), neo4j_data_path + "/CTE_Nodes.csv"  
)  
print("Exec time --- %s seconds ---" % exec_time)
```

Section 7 Analytical Queries

```
MATCH p=()-[r:MENTIONED_IN|affected_players]->() RETURN p LIMIT 25
```

```
MATCH p=()-[r:MENTIONED_IN|suspected_deceased_players]->() RETURN p LIMIT 25
```

```
MATCH p=()-[r:cte_als_former_players|MENTIONED_IN]->() RETURN p LIMIT 25
```

```
MATCH p=()-[r:suspected_deceased_players|MENTIONED_IN]->() RETURN p LIMIT 25
```

Section 8 Project Files

GitHub Repository

<https://github.com/mona-jandro-camm/dse203>

Project Slides

<https://docs.google.com/presentation/d/161il7paAIHnHalWCMS4KxR9bqgB1HZ0-UiOPB2FHUZg/edit?usp=sharing>

Section 9 Revision History

Version	Date	Name	Description
1.0	12-6	Camm Perera	Initial Document Creation
1.1	12-7	Mona Henry	Data Source 3 integration
1.3	12 -10	Mona Henry	Update Knowledge Graph section