Student Name: Megha Nagabhushan

ClassID: 15

We are supposed to build a knowledge graph for the following input (text data).
1. Describe your knowledge about knowledge graph.
2. Why do we want to build such a knowledge graph?
3. What steps are required? Show your own workflow for this task.
4. What are the challenges?
5. Draw a knowledge graph for the given data.

---

CHICAGO (AP) –  Citing high fuel prices, United Airlines said Friday it has increased fares by $6 per round trip on flights to some cities also served by lower-cost carriers. American Airlines, a unit AMR, immediately matched the move, spokesman Tim Wagner said. United, a unit of UAL, said the increase took effect Thursday night and applies to most routes where it competes against discount carriers, such as Chicago to Dallas and Atlanta and Denver to San Francisco, Los Angeles and New York.

---

① Knowledge graphs are graphs of interlinked entities and their attributes. Knowledge graphs have an ontology as its schema.

② a) The knowledge graph are made of interlinked entities. It provides a straightforward approach where users can understand these entities by directly providing key facts. (For example, information about Kansas City)

b) For questions for which answers are not directly available, the users can use efficient exploitation techniques and locate the right portion of the knowledge graph to get the answer.

(For example, what were the 5 coldest years in Kansas City in the last 15 years?)

c) Linking multiple knowledge graphs might help in decision making in large organization.

→ we are improving the learning curve.

→ graph is the most flexible data structure. So, by building a knowledge graph.

querying and information retrieval will be easy.

→ minimizes redundancy.

→ improves decision-making.

③ Step: 1   Natural Language Processing.

Annotation → Lemmatization → POS Tagging → Named Entity Recognition.

CHICAGO ⟶ CHICAGO ⟶ NNP ⟶ LOCATION

-LRB- ⟶ -lrb- ⟶ -LRB-

AP ⟶ AP ⟶ NER ⟶ ORGANIZATION.

-RRB- ⟶ -rrb- ⟶ -RRB-

- ⟶ - ⟶ :

Citing ⟶ cite ⟶ VBG

high ⟶ high ⟶ JJ

fuel ⟶ fuel ⟶ NN

prices ⟶ price ⟶ NNS

` ⟶ `

United → United → NNP → ORGANIZATION

Annotation → Tokenization → $\underline{\text{POS Tagging}}$ → $\underline{\text{Named Entity Recognition}}$

Airlines → Airlines → NNPS → ORGANIZATION

said → say → VBD

Friday → Friday → NNP → DATE

it → it → PRP

has → have → VBZ

increased → increase → VBN

fares → fare → NNS

by → by → IN

$ → $ → MONEY

6 → 6 → CD → MONEY

per → per → IN

round → round → NN

trip → trip → NN

on → on → IN

flights → flight → NNS

to → to → TO

some → some → DT

cities → city → NNS

also → also → RB

served → serve → VBD

by → by → IN

lower-cost → lower-cost → JJ

carriers → carriers → NNS

. → . → .

American → American → NNP. → ORGANIZATION

Airlines → Airlines → NNPS → ORGANIZATION

, → , → ,

a → a → DT

unit → unit → NN

AMR → AMR → NNP

, → , → ,

immediately → immediately → RB

matched → match → VBD

the → the → DT

move → move → NN

, → , → ,

spokesman → spokesman → NN

Tim → Tim → NNP → PERSON

Wagner → Wagner → NNP → PERSON

said → say → VBD

. → . → .

United → United → NNP → ORGANIZATION

, → , → ,

a → a → DT

unit → unit → NN

of → of → IN

UAL → UAL → NNP → ORGANIZATION

, → ,

said → soy → VBD

the → the → DT

increase → increase → NN

took → take → VBD

effect → effect → NN

Thursday → Thursday → NNP → DATE.

night → night → NN → TIME

and → and → CC

applies → apply → VBZ

to → to → TO

most → most → JJS.

routes → route → NNS

where → where → WRB

it → it → PRP

competes → compete → VBZ

against → against → IN

discount → discount → NN

carriers → carries → NNS

, → ,

such → such → JJ

as → as → IN

Chicago → Chicago → NNP → LOCATION.

to $\rightarrow$ to $\rightarrow$ TO

Dallas $\rightarrow$ Dallas $\rightarrow$ NNP $\rightarrow$ LOCATION

and $\rightarrow$ and $\rightarrow$ CC

Atlanta $\rightarrow$ Atlanta $\rightarrow$ NNP $\rightarrow$ LOCATION

and $\rightarrow$ and $\rightarrow$ CC

Denver $\rightarrow$ Denver $\rightarrow$ NNP $\rightarrow$ LOCATION

to $\rightarrow$ to $\rightarrow$ TO

San $\rightarrow$ San $\rightarrow$ NNP $\rightarrow$ LOCATION

Francisco $\rightarrow$ Francisco $\rightarrow$ NNP $\rightarrow$ LOCATION

, $\rightarrow$ , $\rightarrow$ ,

Los $\rightarrow$ Los $\rightarrow$ NNP $\rightarrow$ LOCATION

Angeles $\rightarrow$ Angeles $\rightarrow$ NNP $\rightarrow$ LOCATION

and $\rightarrow$ and $\rightarrow$ CC

New $\rightarrow$ New $\rightarrow$ NNP $\rightarrow$ LOCATION

. $\rightarrow$ . $\rightarrow$ .

Natural language processing refers to the use and ability of systems to process sentences in a natural language such as English, rather than in a specialized artificial computer language such as C++.

It involves the following steps :

① Tokenization : Tokenizing words start into words.

② Lemmatization/Stemming : st

Stemming is the process of reducing words to their stem, base or root form - generally a written word stem.

Lemmatization is the process of grouping together the different inflected forms of a word so they can be analysed as a single term.

c) POS Tagging / Chunking

Parts-of-Speech Tagging. Uses a pre-defined Tagset for this process.

d) Parsing

It is the process of analyzing a sentence by taking each word and determining its structure. from its constituent parts. It uses two components: a parser and a grammar.

e) Name Entity Recognition.

It is a subtask of Named-Entity-Recognition. Information extraction that seeks to locate and classify named entities in text into pre-defined categories such as names of persons, organizations, locations, expressions of times, quantities, monetary values, percentages etc.

f) Co-Reference Resolution

It is the task of finding all expressions that refer to the same entity in a text. It is important step for higher NLP tasks like document summarization, question answering and information extraction.

· Step 2: Information Retrieval and Extraction.
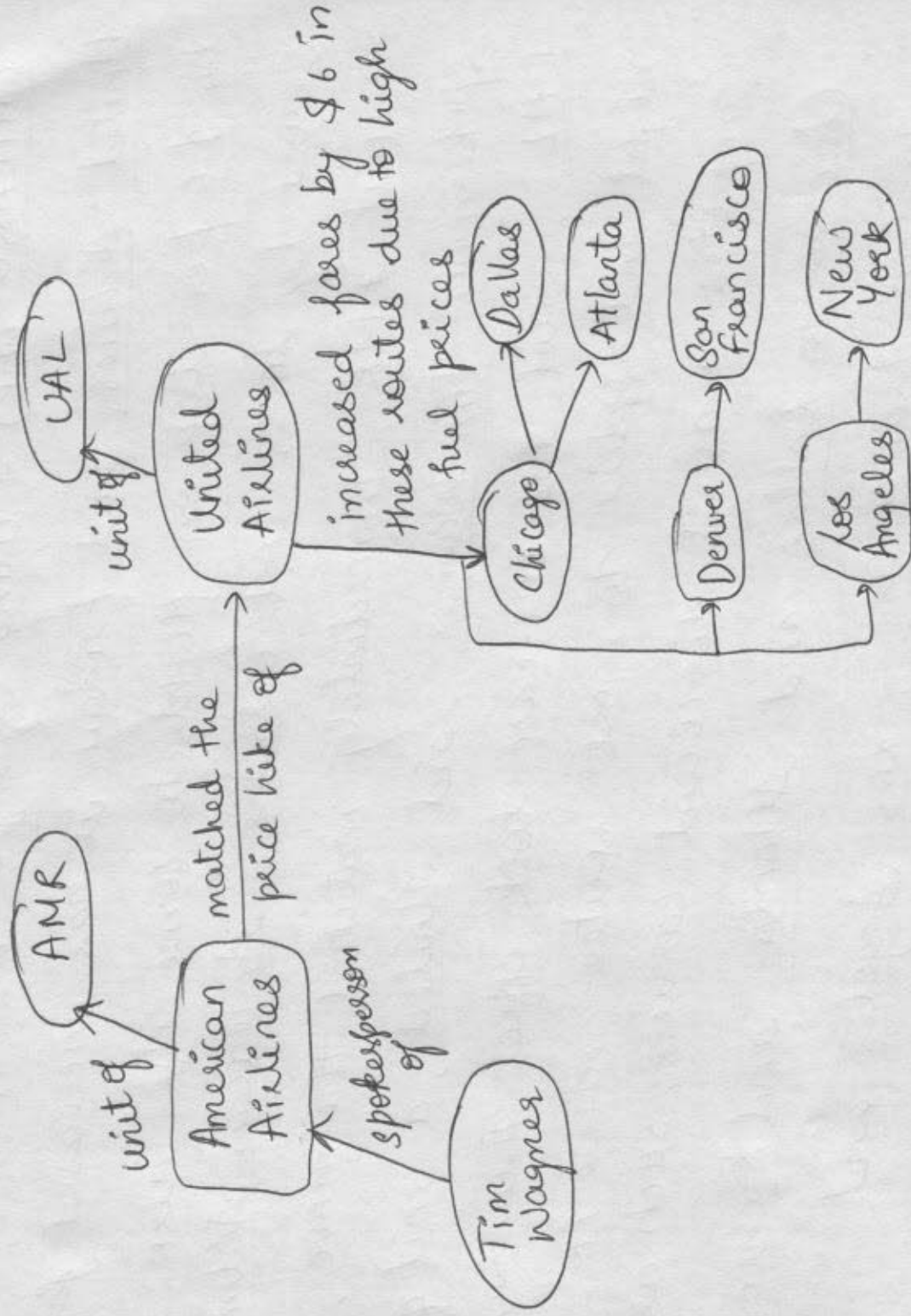
CHICAGO (AP) - Citing high fuel prices, United Airlines said Friday it has increased fares by $6 per trip on flights to some cities also served by lower-cost carriers. American Airlines, a unit AMR, immediately matched the move, spokesman Tim Wagner said. United, a unit of UAL, said the increase took effect Thursday night and applies to most routes where it competes against discount carriers, such as Chicago to Dallas and Atlanta and Denver to San Francisco, Los Angeles and New York.

Step 3: Topic discovery.

| Organization | People | Places |
|---|---|---|
| United Airlines | Tim Wagner. | Chicago |
| American Airlines | | Dallas |
| AMR | | Atlanta |
| UAL | | Denver |
| | | San Francisco |
| | | Los Angeles |
| | | New York |

Step : 4    Knowledge Graph Constuction .



AMR

unit of

American Airlines → matched the → United Airlines

UAL

unit of

Tim Wagner → spokesperson of → American Airlines

United Airlines → increased fares by $6 in these routes due to high fuel prices → Chicago

Chicago → Dallas

Chicago → Atlanta

Denver → San Francisco

Los Angeles → New York

price hike of

4) Challenges of Knowledge Graph.

There is no universal algorithm to generate these knowledge graphs. As a result the knowledge graphs from different sources will adhere to different formats.

Completeness of the knowledge Graph. There can be a lot of missing information.

Correctness of the knowledge Graph. The knowledge Graph may not have the accurate data.

Public knowledge bases may not have the accurate data.

5) Draw knowledge Graph.

P. T. O.