Information Extraction – #1 פרוייקט

בפרוייקט זה תבנו מערכת למענה על שאלות בשפה טבעית בנושא גיאוגרפיה, תוך שימוש בידע שלכם על אונטולוגיות, HTML, SPARQL ו-Xpath. התרגיל להגשה עד ה-01.06, וכמו כל תרגילי הבית, יש להגישו בזוגות. תרגיל זה מהווה 11% מהציון הסופי בקורס.

תיאור המערכת

על המערכת לדעת לענות על שאלות בשפה האנגלית, כאשר כל השאלות יהיו תמיד מאחת התבניות הבאות:

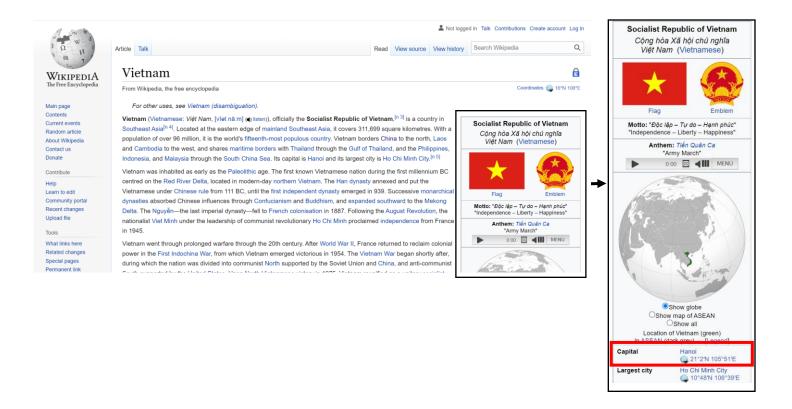
- 1. Who is the president of <country>?
- 2. Who is the prime minister of <country>?
- 3. What is the population of <country>?
- 4. What is the area of <country>?
- 5. What is the form of government in <country>?
- 6. What is the capital of <country>?
- 7. When was the president of <country> born?
- 8. Where was the president of <country> born?
- 9. When was the prime minister of <country> born?
- 10. Where was the prime minister of <country> born?
- 11. Who is <entity>?
- 12. How many <government form1> are also <government form2>?
- 13. List all countries whose capital name contains the string <str>
- 14. How many presidents were born in <country>?

<u>בנוסף,</u> עליכם להגדיר תבנית נוספת לבחירתכם שמסתמכת על המידע הקיים במערכת, ולאפשר למשתמש לשאול שאלות בתבנית זאת. על התבנית להכיל לפחות משתנה אחד.

השאלות יכולות להכיל התייחסויות לשלושה סוגי משתנים:

- Substring: תת מחרוזת כלשהי, מופיעה בשאלה 13 בלבד.
- - של הישות. Wikipedia Infobox: בל יחס הוא שדה ב-Relation

למשל התשובה לשאלה: ?What is the capital of Vietnam, כאשר המידע על היחס מגיע מהשדה המסומן ב-Infobox.



איסוף המידע ובניית האונטולוגיה

עליכם לאסוף מידע על המדינות המופיעות בעמוד הזה:

https://en.wikipedia.org/wiki/List of countries by population (United Nations)

שימו לב שעליכם לחלץ מידע לא רק מה-infobox בעמודי המדינות, אלא גם מהעמודים של ראשי/ות הממשלה והנשיאות/ים. השתמשו בידע שלכם על Crawlers ו-Xpath כדי לעבור בצורה אוטומטית על הדפים הרלוונטים ולחלץ משם את המידע הדרוש.

את האונטולוגיה יש לשמור בקובץ בשם ontology.nt את האונטולוגיה יש לשמור

מענה על שאלות בשפה טבעית

בהנתן שאלה באנגלית, על התוכנית לתרגם את השאלה לשאילתת SPARQL שתרוץ מעל האונטולוגיה שבניתם ותחזיר את התשובה. התשובה לא צריכה להיות "תשובה מלאה", אלא רק להכיל את הערך הנדרש.

- למשל עבור השאלה: ?Where was Justin Trudeau born התשובה תהיה:Canada אין צורך לציין מידע מעבר לשם המדינה. אם שם המדינה לא מצויין, אין צורך לאסוף מידע אחר.
- עבור השאלה: ?Who is Pedro Castillo ושם המדינה בה הוא מחזיק בתפקידו: President) ושם המדינה בה הוא מחזיק בתפקידו: President of Peru

שימו לב ששאלות מתבנית זאת ישאלו תמיד על ראשי/ות ממשלה או נשיאות/ים.

בשאלות מסויימות יכולה להיות יותר מתשובה אחת, כמו למשל: What is the form of government in Argentina? בשאלות מסויימות יכולה להיות יותר מתשובה אחת, כמו למשל: במקרה כזה נציג את כל התשובות מופרדות פסיקים (רווח אחרי כל פסיק), וממויינות בסדר לקסיקוגרפי: Federal republic, Presidential system, Republic

הרצת הקוד

- על הקוד להיות כתוב בפייתון 3 ולרוץ באופן תקין בנובה.
- פותרוץ משורת הפקודה באופן הבא: geo_qa.py •
 - python3 geo_qa.py create o במצב create התכנית תייצר את הק

במצב create התכנית תייצר את הקובץ ontology.nt שיביל את האונטולוגיה שבניתם ותסיים לרוץ.

- י "python3 geo_qa.py question "<question>" מפסבים קמפר שאלה ותסיים לרוץ. במצב question התכנית תקבל שאלה בשפה טבעית, תדפיס למסך את התשובה לשאלה ותסיים לרוץ. השאלה ניתנת כמחרוזת אחת, כלומר מועברת בשורת הפקודה כמחרוזת שמתחילה במרכאות ומסתיימת במרכאות.
 - על התכנית להסתיים לאחר הרצת הפקודה (create). אין להשאיר את התכנית רצה.

תיאור הפרוייקט

עליכם להגיש קובץ נוסף בשם project.pdf שיכיל את הפרטים הבאים

- שמות ומספרי הת.ז של המגישים
- תיאור של הקוד שבונה את האונטולוגיה flow וחלקים חשובים.
- תיאור של השאלה שהוספתם למערכת ודוגמאות לתשובות אפשריות.
- תיאור של שלושה מקרי קצה שהתמודדתם איתם באיסוף המידע. כמו שראינו בתרגול, יתכנו מקרים ספציפיים שידרשו טיפול מיוחד מקרים בהם הייתם צריכים לכתוב שאילתות נוספות כדי להתמודד עם חלק קטן ב-xpath שמופיע בפורמט שונה מהשאר. תארו 3 מקרים כאלו שנתקלתם בהם במהלך העבודה הסבירו אילו חיפושי מיוחדים הייתם צריכים להוסיף ואיך המבנה של המקרה הזה היה שונה ממקרים אחרים.

הוראות הגשה

עליכם להגיש קובץ zip בשם hw1 <id1> <id2>.zip שיכיל את הקבצים הבאים:

- geo_qa.py .1 הקובץ שמכיל את התכנית שבונה את האונטולוגיה ועונה על השאלות
 - ontology.nt קובץ אונטולוגיה בנוי
 - project.pdf .3 תיאור הפרוייקט
- 4. requirements.txt הספריות הנדרשות להרצת הפרוייקט, ראו פירוט ב"עבודה עם ספריות חיצוניות".

אין בעיה לפצל את הקוד למספר קבצים ולהוסיף קבצי עזר כל עוד הקוד עובד כמצופה. במקרה כזה יש להגיש את כל הקבצים הרלוונטים. קבצים שאינם zip לא יבדקו.

בדיקת הפרוייקט

ציון הפרוייקט מחושב באופן הבא:

- דיהות גלויות 74%
- 16% בדיקות נסתרות
- 10% תיאור הקוד, מקרי הקצה, והשאלה שהוספתם למערכת

ייבדקו 45 שאלות בשפה טבעית. 37 מהשאלות זמינות לכם במודל ותוכלו לבדוק את הקוד שלכם עליהן. 8 השאלות הנוספות נסתרות.

הפרוייקט יבדק באופן אוטומטי בנובה, לכן אנחנו ממליצים להקפיד שהקוד רץ ללא שגיאות ועומד במבנה התשובות הנדרש. תשובות בפורמט אחר, או שונות מהתשובות המצופות אפילו בתווים בודדים ייחשבו כתשובה שגויה.

המידע בויקיפדיה משתנה עם הזמן ויכולים להיות שינויים בערכים שרלוונטים לפרוייקט. אם אתם חושבים שהתשובה באחת השאלות הגלויות השתנתה, בבקשה תכתבו לנו בפורום כדי שנוכל לעדכן. בכל מקרה, הקוד שלכם יבדק מול גרסת האונטולוגיה שהגשתם, כך שאם סיימתם לפני הזמן תוכלו להגיש את הפרוייקט ללא חשש מעדכונים בויקיפדיה.

דוגמאות הרצה

מצורפות מספר דוגמאות הרצה:

nova:~> python3 geo_qa.py question "Who is the president of Portugal?" Marcelo Rebelo de Sousa

nin3.7 install flask

nova:~> python3 geo_qa.py question "What is the form of government in Sweden?" Constitutional monarchy, Parliamentary system, Unitary state

nova:~> python3 geo_qa.py question "List all countries whose capital name contains the string hi" Bhutan, India, Moldova, Sint Maarten, United States

עבודה עם ספריות חיצוניות

בדי לאפשר עבודה עם ספריות חיצניות בנובה, יש לעבוד עם סביבה וירטואלית. אפשר למצוא הוראות להרמת סביבה באן.

צור סביבה וירטואלית (כאן נשתמש בפייתון 3.7 ופיפ 3.7, ניתן לשנות זאת אם רוצים) (זה צריך להתבצע רק נשאין סביבה)

virtualenv --prompt=<env-prefix> --python=python3.7 <env-path>

virtualenv --prompt=<my-env> --python=python3.7 .env

virtualenv --prompt=<my-env> --python=python3.7 .env

reduction for a control of the c

יש לצרף להגשה קובץ בשם requirements.txt, שיכיל את שמות כל הספריות הנדרשות להרצת הפרוייקט, כל ספריה בשורה נפרדת. אפשר לקרוא עוד ולראות דוגמא <u>כאן</u>. אפשר לייצר את הקובץ בעזרת הפקודה freeze:

- היכנסו לסביבת הפייתון בה כתבתם את התרגיל
 - הריצו את השורות הבאות:

```
from pip._internal.operations import freeze
print('\n'.join(freeze.freeze()))
```

• תוכן הקובץ יהיה הפלט של ההדפסה

בהצלחה!