

תרגיל תכנותי #1 – Information Extraction

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

תיאור המערכת

על המערכת לדעת לענות על שאלות מהסוגים המפורטים למטה. כל השאלות יהיו באנגלית ויהיו תמיד מאחת התבניות הבאות.

שאלות על ישויות ספציפיות:

1. Who **directed** <film>?
2. Who **produced** <film>?
3. Is <film> **based on** a book?
4. When was <film> **released**?
5. How long is <film>?
6. Who **starred in** <film>?
7. Did <person> **star in** <film>?
8. When was <person> **born**?
9. What is the **occupation** of <person>?

שאלות כלליות:

1. How many films are **based on** books?
2. How many films **starring** <person> won an academy award?
3. How many <occupation1> are also <occupation2>?

שאלה לבחירתכם:

עליכם להוסיף שאלה נוספת לבחירתכם, שמסתמכת על המידע הקיים במערכת.

השאלות יכולות להכיל התייחסויות לשני סוגי משתנים:

- **Entity**: ישות שיש לה ערך בויקיפדיה. לדוגמא לישות Emma Watson יש את הדף https://en.wikipedia.org/wiki/Emma_Watson. שם הישות יהיה זהה לשמה ב-URL של דף הויקיפדיה שלה עם רווח במקום קו תחתון.
- **Relation**: כל יחס הוא שדה ב-Wikipedia Infobox של הישות.
למשל התשובה לשאלה: Who **directed** The Great Gatsby (2013 film)? תהיה Baz Luhrmann, כאשר המידע על היחס מגיע מהשדה המסומן ב-Infobox.



WIKIPEDIA
The Free Encyclopedia

- 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
- Page information
- Cite this page
- Wikidata item
- Print/export
- Download as PDF
- Printable version

<https://en.wikipedia.org/wiki/Special:SpecialPages>

Not logged in

Talk

Contributions

Create account

Log in

Article

Talk

Read

Edit

View history

The Great Gatsby (2013 film)

From Wikipedia, the free encyclopedia

The Great Gatsby is a 2013 romantic drama film based on F. Scott Fitzgerald's 1925 novel of the same name. The film was co-written and directed by Baz Luhrmann and stars Leonardo DiCaprio as the eponymous Jay Gatsby, with Tobey Maguire, Carey Mulligan, Joel Edgerton, Isla Fisher, Jason Clarke, Elizabeth Debicki and Jack Thompson.^[4] Jay-Z served as executive producer. Production began in 2011 and took place in Australia, with a \$105 million net production budget. The film follows the life and times of millionaire Jay Gatsby (DiCaprio) and his neighbor Nick Carraway (Maguire), who recounts his encounter with Gatsby at the height of the Roaring Twenties on Long Island.

The film was highly polarizing among critics; it received alternating praise and criticism for its acting performances, soundtrack, visual style, and direction. Audiences responded more positively^[5] and Fitzgerald's granddaughter praised the film, stating "Scott would have been proud."^[6] As of 2017, it is Luhrmann's highest-grossing film, grossing over \$353 million worldwide.^[7] At the 86th Academy Awards, the film won in both of its nominated categories: *Best Production Design* and *Best Costume Design*.

Contents [hide]

1 Plot

2 Cast

3 Production

3.1 Development

3.2 Casting

3.3 Screenplay

3.4 Filming

3.4.1 Sets

3.4.2 Costumes

4 Release and marketing

5 Soundtrack

The Great Gatsby



Theatrical release poster

Directed by

Baz Luhrmann

Produced by

Baz Luhrmann

Catherine Knapman

Douglas Wick

Lucy Fisher

Catherine Martin

Screenplay by

Baz Luhrmann

Craig Pearce

איסוף המידע ובניית האונטולוגיה

עליכם לאסוף מידע על הסרטים המופיעים בעמוד הזה:

https://en.wikipedia.org/wiki/List_of_Academy_Award-winning_films

שימו לב שעליכם לחלץ מידע לא רק מה-infobox בעמודי הסרטים, אלא גם בעמודים של המפיקים/ות, במאים/ות והשחקנים/שחקניות. השתמשו בידע שלכם על SPARQL ו-Xpath כדי לעבור בצורה אוטומטית על הדפים הרלוונטים ולחלץ משם את המידע הדרוש.

העמוד מכיל מספר רב של סרטים, ואנחנו נבדוק את איסוף המידע רק על סרטים ששנת הזכיה שלהם (לפי הטבלה הראשית) היא 2010 (כולל).

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

מענה על שאלות בשפה טבעית

על התוכנית לדעת להתמודד עם שאלות באנגלית על גבי האונטולוגיה. בהנתן שאלה באנגלית (מאחד מ-12 המבנים למעלה) על התוכנית לתרגם את השאלה לשאלת SPARQL שתורץ מעל האונטולוגיה שבניתם ותחזיר את התשובה. התשובה לא צריכה להיות "תשובה מלאה", אלא רק להכיל את הערך הנדרש. למשל עבור השאלה:

Is [Little Women \(2019 film\)](#) based on a book?

התשובה תהיה: Yes
אין צורך לציין את שם הספר.

בשאלות מסויימות יכולה להיות יותר מתשובה אחת, כמו למשל:

Who [produced The Great Gatsby \(2013 film\)](#)?

במקרה כזה נציג את כל התשובות מופרדות פסיקים (רווח אחרי כל פסיק), וממיינות בסדר לקסיקוגרפי:

[Baz Luhrmann](#), [Catherine Knapman](#), [Catherine Martin \(designer\)](#), [Douglas Wick](#), [Lucy Fisher](#)

הרצת הקוד

- על הקוד להיות כתוב בפייתון 3 ולרוץ באופן תקין בנובה. יש לפרט בקובץ requirements.txt כל ספריה חיצונית שנעשה בה שימוש (שמות כל החבילות מופרדות ע"י ירידת שורה)
 - התכנית תיקרא film_qa.py ותרץ משורת הפקודה באופן הבא:
 - `python film_qa.py create`במצב create התכנית תייצר את הקובץ ontology.nt שיכיל את האונטולוגיה שבניתם ותסיים לרוץ.
 - `python film_qa.py question "<question>"`
- 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 או question). אין להשאיר את התכנית רצה.

תיאור הפרוייקט

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

- שמות ומספרי התז של המגשים
- תיאור של הקוד שבונה את האונטולוגיה, איך הוא בנוי ומה עשיתם
- תיאור של השאלה שהוספתם למערכת ודוגמאות לתשובות אפשריות
- תיאור של שלושה מקרי קצה שהתמודדתם איתם באיסוף המידע. הסבירו אילו חיפוישי xpath 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:

1. film_qa.py - הקובץ שמכיל את התכנית שבונה את האונטולוגיה ועונה על השאלות
2. ontology.nt – קובץ אונטולוגיה בנוי
3. project.pdf – תיאור הפרוייקט
4. requirements.txt – קובץ txt שמכיל את שמות כל החבילות החיצוניות שנחוצות להרצת הפרוייקט, מופרדות ע"י ירידת שורה.

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

בדיקת הפרוייקט

ייבדקו 20 שאלות בשפה טבעית. 12 מהשאלות זמינות לכם במודל כך שמובטח שאם התכנית שלכם עונה עליהן בצורה תקינה תקבלו ציון עובר. 8 השאלות הנוספות נסתרות.

הפרוייקט יבדק באופן אוטומטי לחלוטין בנובה, לכן אנחנו ממליצים להקפיד שהקוד רץ ללא שגיאות ועומד במבנה התשובות שנדרש. 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.

בהצלחה!