

**אופציה 2 לתרגיל בית מס' 5
ביצוע ניסוי מעבדה לחקר ה-RAG**

Dr. Yoram Segal

כל הזכויות שמורות

26 בנובמבר 2025

תוכן העניינים

1	מבוא
2	הקדמה
2	מטרת הניסוי
2	מבנה הניסוי: 3 תת-ניסויים
2	ניסוי 1: "הבעיה" -- חלון הקשר וגודל הקשר
2	ניסוי 2: "האשליה" -- גישה וחוסר רלוונטיות
3	ניסוי 3: "הפתרון" -- שימוש ב-RAG
	шиפורים ושידרגים לטיב ניתוח הניסויים -- נקודות להתייחסות בכל אחד מהניסויים
3	3 הוראות לסטודנט
3	שלב 1 -- ניסוי "קלאסי" ללא RAG
3	שלב 2 -- הוספה רعش
4	שלב 3 -- ביצוע RAG
4	דגשים לחלק המעשי
4	Ollama
4	LangChain
5	טבלת תכנון הניסוי
6	הצגת תוצאות הניסוי

1 מבוא

מסמך זה מציג תכנון ניסוי ללימוד והדגמת השפעת RAG (Retrieval-Augmented Generation) על בעיות חלון ההקשר.

1.1 הקדמה

מטרלה זו עוסקת בניתוח ולימוד של RAG. שורת הניסויים המוצעת להלן מהוות מסגרת ריעיונית כללית, ואתם מוזמנים לפרש, לפתח ולחזור את הנושאים בכל דרך שתמצאו לנו. עבור כל ניסוי עלייכם להגדיר שאלות מחקר, לבצע את הניסויים ולנתח את הממצאים, רצוי תוך הצגת ניתוח סטטיסטי ויזואלי (באמצעות גרפים או טבלאות). מומלץ לחזור על כל ניסוי מספר פעמים כדי להבטיח תוקף סטטיסטי לתוצאות.
שים לב: המסקנות שלכם אינן חיבות לחופוף לחומר שהוצע בשיעור. אתם רשאים להגיע לתובנות עצמאיות, ובלבך שתתמקו אותן היטב; במקרים אלו מומלץ להיעזר בסימוכין חיצוניים ולהציג הסבר לעירומים שנגילותם. קחו את הניסויים למקומות שימושיים אתם ולכיוון החקירה האישי שלכם נועדו לשמש כ'סיעור מוחות' ואין בגדיר הגדרות סגורות.

1.2 מטרת הניסוי

להמבחן לסטודנט, בקוד Python קצר, מדוע נוצרות בעיות "איבוד באמצע" (lost in the middle) והזיות במודול שפה כאשר הקשר ארוך מדי או בלתי רלוונטי -- וכייד RAG פותר זאת, על-ידי אחזור ממוקד.

2 מבנה הניסוי: 3 תת-ניסויים

2.1 ניסוי 1: "הבעיה" -- חלון הקשר וגודל הקשר

-- **צירת תוכן סינטטי:** שורת עובדות פשוטות, לדוגמה: ["Paris is the capital of France.", "Berlin is the capital of Germany."], ...]

-- **שיז שאלת:** כל פעם מוסיפים עובדות חדשות ל"קשר" ובזקדים האם המודל מצליח לענות על שאלת עובדה שמוסתרת באמצע/בסוף/בהתחלת ("מהי בירת גרמניה?").

-- **מדידה:** אחוז הדיק/הצלחת שליפת המידע לפי מיקום העובדה והאורך.

2.2 ניסוי 2: "האשליה" -- גלישה וחוסר רלוונטיות

-- **הוספת "רעש":** לרשימה העובדות נוסיף משפטים לא רלוונטיים או הרבה "פילרים".

-- **בדיקה:** שוב נבקש את המודל לענות -- האם הצלחה יורדת כי הקשר "רווי" ומבלב?

-- מדידה: איקוח התשובה, הופעת "זהיה" או דיק ירוד.

2.3 ניסוי 3: "הפתרון" -- שימוש ב-RAG

-- הטמעה באמצעות Embeddings: כל עובדה מקבלת embedding (על-פי-nomic) .LangChain או Sentence Transformers embed-text Ollama

-- **בנייה DB**: שימוש ב-Chroma/FAISS וcdcמה -- לאחזר אנלוגי בסגנון RAG.

-- **שאלתה**: שואלים את אותה שאלה, אך הקוד שלפָר רק את העבודה הци רלוונטי מה-DB.

-- מדידה: האם עכשו גם בשדה "רועל", RAG מוצא (כמעט תמיד) את העבודה הרלוונטי?

2.4 שיפורים ושידורגים לטיב ניתוח הניסויים -- נקודות להתייחסות בכל אחד מהניסויים

-- **Reranking**: דירוג משופר בעזרת מודל LLM.

-- **הוספת הקשר**: הוספה הקשר לפסקאות או הגדרה מושכלת.

-- **טעויות RAG נפוצות**: החסارة בעובדה, מסמך גדול מדי, מידע "נתקע" בהקשר.

3 דף הוראות לסטודנט

3.1 שלב 1 -- ניסוי "קלאס" ללא RAG

1. צור 02--03 "מסמכים בעובדה" סינטטיים.

2. שמור את כולם כחיבור-טקסט ארוך אחד.

3. הגש למודל בעזרת ollama.generate או chat ומצג שאלות כמודגמים לעיל.

4. בדוק: האם המידע באמצע או בסוף נשאר זמין ונשלף?

3.2 שלב 2 -- הוספת רعش

1. הוסף טקסטים רנדומליים באורך רב.

2. נסהשוב לשאול. בדוק: איך הביצועים?

RAG שלב 3 -- ביצוע

1. עבור כל עובדה: הפק embedding (Ollama/LangChain) ובודקה: האם קטע רלוונטי, ורק אותו תן למודל.
2. בנה Chroma/FAISS Vector DB (וכדומה).
3. בעת שאלתה: שאל את רשות הוקטורים, החזר קטע רלוונטי, ורק אותו תן למודל השפה קשר.
4. מודוד: האם איקוט התשובה משתפרת? האם המודל פחות מזין?

4 דגשים לחלק המעשי

Ollama 4.1

השתמש בפקודות כמו `ollama.generate`, `ollama.embeddings` ובמידת הצורך `Chroma` Python DB.

LangChain 4.2

הפעיל את ההتمמשקות עם Chain, EmbeddingModels, Chroma ו-Retriever המובנים.
הערה חשובה: בכל צעדי -- מודיעו במפורש את מהירות המענה, אורך הקשר ו-דיקות התשובה.

5 טבלת תכנון הניסוי

יעסינ-תת	סינון המ	סילואש המ	האוושה/הייפיצ
1 יוסין	תודבען זורא זמסם	הdboע לע להאל X סוקימבר	דרוי קויד פוס/עכמאנב
2 יוסין	+ שער אלם זמסם תודבען	הdboעה לע להאל	טיזה, קויד תדיiri
(RAG)	Embeddings ירוטקו רוזחא	רוזחא -> להאל החילש הדבועה, LLM	ההובג החלצה רתוי רהמ

טבלה 1: טבלת תכנון תתי-הניסויים

6 הציגת תוצאות הניסוי

על הסטודנט לתקן ולהשוו באיזה אופן משכנע להציג את תוצאות הניסוי, החקיר של הניסוי, והמסקנות מהניסוי. מומלץ לתקף את התוצאות בגרפים לפי שיקול דעת הסטודנט.

הערה: האמור במסמך זה מיועד לנשים ולגברים כאחד, והשימוש בלשון זכר הוא מטעמי נוחות בלבד.