

حَسْنَةٌ لِرَحْمَةِ رَبِّنَا



بازیابی اطلاعات

پروژه: ارزیابی مدل‌های فضای برداری و مدل‌های احتمالی زبانی

آرمان آزادی

۱۴۰۴ زمستان

فهرست مطالب

چکیده ۴

۵	مقدمه	.۱
۵	داده‌ها و پیش‌پردازش	.۲
۸	روش‌شناسی: مدل‌های بازیابی	.۳
۱۲	ارزیابی	.۴
۱۴	نتایج تجربی	.۵
۱۹	بحث و تحلیل	.۶
۲۰	نتیجه گیری	

چکیده

این گزارش به پیاده‌سازی، تنظیم و ارزیابی یک سیستم بازیابی اطلاعات می‌پردازد که برای بازیابی اسناد مرتبط با پرسش‌های کاربر طراحی شده است. این سیستم تحت محدودیت‌های خاصی ساخته شده است، از جمله محدودیت اندازه واژگان به حداقل ۵۰۰ کلمه و ممنوعیت استفاده از کتابخانه‌های خارجی. سه مدل بازیابی متمایز پیاده‌سازی شدند: مدل فضای برداری (BM25)، مدل زبانی تک‌کلمه‌ای (Unigram) با هموارسازی (Dirichlet Smoothing)، و مدل زبانی دو‌کلمه‌ای (Bigram) با درونیابی (Interpolation). مدل‌ها بر روی یک مجموعه داده اعتبارسنجی (Validation Set) برای بهینه‌سازی ابرپارامترها (λ , k_1 , b , μ) تنظیم شدند و سپس بر روی یک مجموعه داده آزمایشی (Test Set) با استفاده از معیارهای دقت در ۵ نتیجه اول ($P@5$)، میانگین رتبه معکوس (MRR) و میانگین دقت متوسط (MAP) ارزیابی شدند. نتایج نشان می‌دهد که مدل BM25 و مدل Bigram درونیابی‌شده عملکردی قابل مقایسه داشتند، و مدل BM25 با دستیابی به بالاترین میانگین دقت متوسط (MAP) در مجموعه آزمایشی، بهترین عملکرد را داشت.

۱. مقدمه

هدف اصلی این پروژه ساخت یک موتور جستجو بود که قادر باشد مجموعه‌ای از اسناد را بر اساس ارتباط آن‌ها با پرسش‌های خاص زبان طبیعی رتبه‌بندی کند.

دامنه این پروژه شامل موارد زیر بود:

۱. پردازش داده‌ها: دریافت مجموعه‌ای از اسناد متنی و پرسش‌های کاربران.
۲. نمایه‌سازی (Indexing): ساخت نمایه‌های معکوس برای امکان بازیابی کارآمد.
۳. مدل‌سازی: پیاده‌سازی سه الگوریتم رتبه‌بندی متفاوت (BM25، Unigram LM، Bigram LM).
۴. ارزیابی: پیاده‌سازی دستی معیارهای استاندارد IR برای سنجش کمی عملکرد.

یک محدودیت کلیدی در این پیاده‌سازی، محدودیت واژگان بود. سیستم محدود به حداقل ۵۰۰ کلمه پرترکار موجود در مجموعه آموزشی بود.

۲. داده‌ها و پیش‌پردازش

آمار مجموعه داده‌ها

مجموعه داده به سه بخش مجزا تقسیم شد تا آموزش و ارزیابی دقیق تضمین شود:

- **مجموعه آموزشی (Training Set)**: ۷۵۰ سند. صرفاً برای ساخت واژگان و تخمین احتمال پس‌زمینه (Background Probability) استفاده شد.
- **مجموعه اعتبارسنجی (Validation Set)**: ۱۲۴۵ سند و ۲۰ پرسش. برای تنظیم ابرپارامترها استفاده شد.
- **مجموعه آزمایشی (Test Set)**: ۱۲۴۵ سند و ۳ پرسش. برای ارزیابی نهایی استفاده شد.

```
Loading Data:  
Train Docs: 750  
Val Docs: 1245, Val Queries: 20  
Test Docs: 1245, Test Queries: 30
```

شکل ۱

خط لوله پردازش متن

داده‌های متنی خام، بدون ساختار هستند. برای آماده‌سازی داده‌ها جهت نمایه‌سازی، یک خط لوله پیش‌پردازش یکسان بر روی تمام اسناد و پرسش‌ها اعمال شد:

۱. **نرمال‌سازی**: تمام متن‌ها به حروف کوچک تبدیل شدند تا حساسیت به بزرگی و کوچکی حروف از بین برود.
۲. **پاک‌سازی**: علائم نگارشی و کاراکترهای خاص با استفاده از عبارات باقاعده (Regex) حذف شدند.
۳. **توکن‌سازی (Tokenization)**: متن بر اساس فضای خالی (Whitespace) به کلمات جداگانه تقسیم شد.
۴. **حذف کلمات توقف (Stopwords)**: لیستی از کلمات رایج انگلیسی (مانند "the", "is", "at") فیلتر شدند تا نویز کاهش یابد و حجم نمایه کم شود.

```
8 def preprocess(text):
9     """
10    Lowercases, removes punctuation, tokenizes, and removes stopwords.
11    """
12    if not isinstance(text, str): return []
13    text = text.lower()
14    # Remove punctuation
15    text = re.sub(r'[^\w\s]', ' ', text)
16    tokens = text.split()
17    return [t for t in tokens if t not in STOPWORDS]
18
```

شکل ۲

محدودیت واژگان

پس از پیش‌پردازش ۷۵۰ سند آموزشی، فرکانس کلمات محاسبه شد و تنها ۵۰۰ کلمه پر تکرار حفظ شدند.

- تاثیر: هر کلمه‌ای در اسناد اعتبارسنجی یا آزمایشی (یا پرسش‌ها) که در این لیست ۵۰۰ تایی نبود، در حین نمایه‌سازی و بازیابی حذف شد. این امر به عنوان یک فیلتر شدید انتخاب ویژگی عمل می‌کند اما ساختار معنایی اصلی متن را حفظ می‌کند.

```

19 def build_vocab(documents, top_k=500):
20     """
21         Builds vocabulary from TRAINING documents using top K frequent words.
22         Constraint: Length of TF-IDF vector must be at least 500.
23     """
24     print(f"Building vocabulary from {len(documents)} training documents:")
25     all_tokens = []
26     for doc_text in documents.values():
27         all_tokens.extend(preprocess(doc_text))
28
29     # Count frequencies
30     counts = Counter(all_tokens)
31
32     # Select top K (K must be at least 500)
33     most_common = counts.most_common(top_k)
34     vocab = {term for term, freq in most_common}
35
36     return vocab, counts
37
38 # Build Vocab from TRAINING data:
39 VOCAB, TRAIN_COUNTS = build_vocab(train_passages, top_k=500)
40 print(f"Vocabulary size: {len(VOCAB)}")
41
42 def filter_query(query_text):
43     """
44         Preprocesses query and filters words not in VOCAB.
45         Constraint: New words should not be considered.
46     """
47     tokens = preprocess(query_text)
48     return [t for t in tokens if t in VOCAB]

```

شکل ۳

۳. روش‌شناسی: مدل‌های بازیابی

مدل BM25 (Best Matching 25)

BM25 یک تابع رتبه‌بندی مبتنی بر چارچوب احتمالی ارتباط (Probabilistic Relevance Framework) است. این مدل با معرفی فرکانس کلمه (Term Frequency Saturation) و نرمال‌سازی طول سند، نسبت به TF-IDF استاندارد بهبود یافته است.

فرمول:

امتیاز یک سند D برای یک پرسش Q به صورت زیر محاسبه می‌شود:

$$\text{Score}(D, Q) = \sum_{q \in Q} IDF(q) \cdot \frac{f(q, D) \cdot (k_1 + 1)}{f(q, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{\text{avgdl}})}$$

فرمول ۱

که در آن:

- $f(q, D)$ فرکانس کلمه q در سند D است.
- $|D|$ طول سند D است.
- Avgdl میانگین طول اسناد در کل مجموعه است.
- K_1 پارامتری است که اشباع فرکانس کلمه را کنترل می‌کند.
- b پارامتری است که میزان نرمال‌سازی طول را کنترل می‌کند.

پیاده‌سازی:

مؤلفه IDF با استفاده از نسخه احتمالی محاسبه شد:

$$IDF(q) = \log \frac{N - n(q) + 0.5}{n(q) + 0.5} + 1$$

فرمول ۲

```
35     # Precompute IDF
36     for term in self.vocab:
37         df = doc_freqs.get(term, 0)
38         # Standard Probabilistic IDF
39         self.idf[term] = math.log((self.doc_count - df + 0.5) / (df + 0.5) + 1)
```

شکل ۴

این فرمول به صورت دستی پیاده‌سازی شد. فاز تنظیم بر بهینه‌سازی پارامترهای K_1 و b تمرکز داشت.

مدل زبانی تک کلمه‌ای (Unigram) با هموارسازی دیریشله

مدل سازی زبانی برای بازیابی اطلاعات، اسناد را بر اساس احتمالی که سند D بتواند پرسش Q را تولید کند، رتبه‌بندی می‌کند. مدل تک کلمه‌ای فرض استقلال کلمات (Bag-of-Words) را در نظر می‌گیرد.

هموارسازی (Smoothing)

برای حل مشکل احتمال صفر (زمانی که یک کلمه پرسش در سند وجود ندارد)، از هموارسازی دیریشله با استفاده از پارامتر μ استفاده شد. این تکنیک، مدل زبانی سند را با یک مدل پس‌زمینه مجموعه (بر اساس داده‌های آموزشی) ترکیب می‌کند.

فرمول:

$$P(w|D) = \frac{c(w, D) + \mu P(w|C)}{|D| + \mu}$$

فرمول ۳

که در آن:

- تعداد تکرار کلمه w در سند D است. $c(w,D)$
- احتمال کلمه w در کل مجموعه آموزشی است. $P(w|C)$
- پارامتر هموارسازی است. μ

```

34     for doc_id in candidate_docs:
35         doc_len = self.doc_lengths[doc_id]
36         log_prob = 0.0
37
38         for term in tokens:
39             # Find TF in current doc (naive search in list)
40             tf = 0
41             for d, c in self.inverted_index[term]:
42                 if d == doc_id:
43                     tf = c
44                     break
45
46             # Dirichlet Smoothing
47             p_wc = self.p_C.get(term, 1e-10)
48             numerator = tf + (self.mu * p_wc)
49             denominator = doc_len + self.mu
50
51             prob = numerator / denominator
52             log_prob += math.log(prob)
53
54         scores[doc_id] = log_prob

```

شکل ۵

مدل زبانی دوکلمه‌ای (Bigram) با درونیابی

مدل دوکلمه‌ای با در نظر گرفتن محلی کلمات (جفت‌های کلمات مجاور)، رویکرد تک‌کلمه‌ای را گسترش می‌دهد. این کار عبارات و ساختار محلی را ضبط می‌کند.

فرمول: از آنجا که پراکندگی داده‌ها زیاد است (به‌ویژه با واژگان ۵۰۰ کلمه‌ای)، یک مدل دوکلمه‌ای خالص اغلب غیرقابل اعتماد است. ما یک مدل درونیابی‌شده (Interpolated) را پیاده‌سازی کردیم:

$$P(q_i|q_{i-1}, D) = \lambda P_{Uni}(q_i|D) + (1 - \lambda) P_{Bi}(q_i|q_{i-1}, D)$$

فرمول ۴

که در آن:

- احتمال هموارشده تک‌کلمه‌ای است. P_{Uni}

- احتمال تجربی ظاهر شدن دوکلمه (Bigram) در سند است. P_{Bi}

- وزن درونیابی است (بین ۰.۰ تا ۱.۰). λ

```

89     # 1. First term: P(q1 | D) using Smoothed Unigram
90     q1 = tokens[0]
91     tf_q1 = doc_seq.count(q1)
92     p_uni_q1 = (tf_q1 + self.unigram_model.mu * self.unigram_model.p_C.get(q1, 1e-10)) / (doc_len + self.unigram_model.mu)
93     log_prob += math.log(p_uni_q1)
94
95     # 2. Subsequent terms: Bigram P(qi | qi-1, D)
96     for i in range(1, len(tokens)):
97         qi = tokens[i]
98         q_prev = tokens[i-1]
99
100        # Count Bigram C(qi, qi-1, D) and Unigram C(qi-1, D)
101        c_bigram = 0
102        for j in range(len(doc_seq) - 1):
103            if doc_seq[j] == q_prev and doc_seq[j+1] == qi:
104                c_bigram += 1
105
106        c_prev = doc_seq.count(q_prev)
107
108        # Empirical Bigram Probability
109        p_bigram_emp = (c_bigram / c_prev) if c_prev > 0 else 0.0
110
111        # Smoothed Unigram Probability for qi
112        tf_qi = doc_seq.count(qi)
113        p_uni_smooth = (tf_qi + self.unigram_model.mu * self.unigram_model.p_C.get(qi, 1e-10)) / (doc_len + self.unigram_model.mu)
114
115        # Interpolation
116        final_prob = (self.chi * p_bigram_emp) + (self.lam * p_uni_smooth)
117
118        if final_prob <= 0: final_prob = 1e-10
119        log_prob += math.log(final_prob)
120
121    scores[doc_id] = log_prob

```

شکل ۶

۴. ارزیابی

برای اندازه‌گیری عینی کیفیت نتایج جستجو، سه معیار استاندارد از پایه پیاده‌سازی شدند:

دقت در K ($P@K$)

نسبت اسناد بازیابی شده مرتبط در ۵ نتیجه اول را اندازه می‌گیرد.

: فرمول

$$P@K = \frac{\text{Number of Relevant Documents in the Top K}}{K}$$

فرمول ۵

```

9   # 1. Precision at K (P@5)
10  retrieved_k = retrieved_doc_ids[:k]
11  hits_k = sum(1 for doc in retrieved_k if doc in relevant_set)
12  p_at_k = hits_k / k

```

شکل ۷

میانگین رتبه معکوس (MRR):

موقعیت اولین سند مرتبط را ارزیابی می‌کند. این معیار سیستم‌هایی را که اولین پاسخ صحیح را در پایین لیست قرار می‌دهند، جریمه می‌کند.

: فرمول

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{rank}_i}$$

فرمول ۶

```

14  # 2. Mean Reciprocal Rank (MRR)
15  mrr = 0.0
16  for i, doc in enumerate(retrieved_doc_ids):
17      if doc in relevant_set:
18          mrr = 1.0 / (i + 1)
19          break

```

شکل ۸

میانگین دقت متوسط (MAP):

یک معیار واحد از کیفیت را در تمام سطوح فراخوانی (Recall) ارائه می‌دهد. این میانگین نمرات دقت متوسط برای هر پرسش است.

: فرمول

$$MAP = \frac{1}{|Q|} \sum_{q=1}^{|Q|} AP(q)$$

فرمول ۷

```

21 # 3. Average Precision (AP)
22 # Constraint: Denominator = number of relevant documents
23 hits = 0
24 sum_prec = 0.0
25 for i, doc in enumerate(retrieved_doc_ids):
26     if doc in relevant_set:
27         hits += 1
28         sum_prec += hits / (i + 1)
29
30 if len(relevant_set) > 0:
31     ap = sum_prec / len(relevant_set)
32 else:
33     ap = 0.0
34
35 return p_at_k, mrr, ap

```

شکل ۹

نکته محدودیت: در محاسبه AP ، مخرج کسر به درستی تعداد کل اسناد مرتبط در نظر گرفته شد تا حتی زمانی که سیستم موفق به بازیابی تمام اسناد مرتبط نمی‌شود، معیار دقیق باشد.

۵. نتایج تجربی

فاز ۱: تنظیم ابرپارامترها (مجموعه اعتبارسنجی)

یک جستجوی شبکه‌ای (Grid Search) بر روی مجموعه اعتبارسنجی انجام دادیم تا پارامترهای بهینه برای هر مدل را تعیین کنیم.

تنظیم BM25

ترکیبات $\{k_1, b\}$ را برای $k_1 \in \{0.5, 0.75, 1.0\}$ و $b \in \{1.2, 1.5, 2.0\}$ آزمایش کردیم.

```
3 # A. Tune BM25
4 # Constraints: Test at least 3 values.
5 k1_vals = [1.2, 1.5, 2.0]
6 b_vals = [0.5, 0.75, 1.0]
7 best_bm25_map = -1
8 best_bm25_params = (1.5, 0.75)
9
10 for k in k1_vals:
11     for b in b_vals:
12         model = BM25(VOCAB, k1=k, b=b)
13         _, _, map_score = evaluate_model(model, val_questions, val_judg, val_passages)
14         print(f"BM25 (k1={k}, b={b}) -> MAP: {map_score:.4f}")
15         if map_score > best_bm25_map:
16             best_bm25_map = map_score
17             best_bm25_params = (k, b)
18
19 print(f"Best BM25 Params: {best_bm25_params}")
```

شكل ۱۰

```
BM25 (k1=1.2, b=0.5) -> MAP: 0.3199
BM25 (k1=1.2, b=0.75) -> MAP: 0.3210
BM25 (k1=1.2, b=1.0) -> MAP: 0.3158
BM25 (k1=1.5, b=0.5) -> MAP: 0.3163
BM25 (k1=1.5, b=0.75) -> MAP: 0.3211
BM25 (k1=1.5, b=1.0) -> MAP: 0.3177
BM25 (k1=2.0, b=0.5) -> MAP: 0.3151
BM25 (k1=2.0, b=0.75) -> MAP: 0.3175
BM25 (k1=2.0, b=1.0) -> MAP: 0.3103
Best BM25 Params: (1.5, 0.75)
```

شكل ۱۱

مشاهده: مدل نسبتاً پایدار بود، اما مقادیر $k_1=1.5$ و $b=0.75$ بالاترین MAP را به دست آوردند.

تنظیم مدل تک کلمه‌ای (Unigram)

پارامتر هموارسازی دیریشله {500, 1000, 2000} را برای mu تغییر دادیم.

```
21 # B. Tune Unigram LM
22 # Constraint: Find appropriate mu.
23 mu_vals = [500, 1000, 2000]
24 best_uni_map = -1
25 best_mu = 1000
26
27 for mu in mu_vals:
28     model = UnigramLM(VOCAB, TRAIN_COUNTS, mu=mu)
29     _, _, map_score = evaluate_model(model, val_questions, val_judg, val_passages)
30     print(f"Unigram (mu={mu}) -> MAP: {map_score:.4f}")
31     if map_score > best_uni_map:
32         best_uni_map = map_score
33         best_mu = mu
34
35 print(f"Best Unigram Mu: {best_mu}")
```

شکل ۱۲

```
Unigram (mu=500) -> MAP: 0.2934
Unigram (mu=1000) -> MAP: 0.2916
Unigram (mu=2000) -> MAP: 0.2899
Best Unigram Mu: 500
```

شکل ۱۳

مشاهده: مقادیر پایین‌تر mu عملکرد بهتری داشتند، که نشان می‌دهد مدل زبانی خاص سند نسبت به مدل پس‌زمینه مجموعه برای این مجموعه داده خاص قابل اعتمادتر بود و mu=500 انتخاب شد.

تنظیم مدل دو کلمه‌ای (Bigram)

وزن درونیابی lambda را تنظیم کردیم، جایی که lambda نشان‌دهنده وزن مؤلفه Unigram است.

```

37 # C. Tune Bigram LM
38 # Constraint: Find optimal lambda.
39 lambda_vals = [0.1, 0.5, 0.9]
40 best_bi_map = -1
41 best_lam = 0.5
42
43 # Note: We reuse the best Mu found above for the unigram component
44 base_unigram = UnigramLM(VOCAB, TRAIN_COUNTS, mu=best_mu)
45
46 for lam in lambda_vals:
47     model = BigramLM(VOCAB, TRAIN_COUNTS, base_unigram, lam=lam)
48     _, _, map_score = evaluate_model(model, val_questions, val_judg, val_passages)
49     print(f"Bigram (lambda={lam}) -> MAP: {map_score:.4f}")
50     if map_score > best_bi_map:
51         best_bi_map = map_score
52         best_lam = lam
53
54 print(f"Best Bigram Lambda: {best_lam}")
55

```

شکل ۱۴

```

Bigram (lambda=0.1) -> MAP: 0.3091
Bigram (lambda=0.5) -> MAP: 0.3091
Bigram (lambda=0.9) -> MAP: 0.3116
Best Bigram Lambda: 0.9

```

شکل ۱۵

مشاهده: بهترین عملکرد با $\text{lambda}=0.9$ به دست آمد. این نشان می‌دهد که مدل به شدت احتمالات Unigram (وزن ۹۰٪) را بر احتمالات Bigram (وزن ۱۰٪) ترجیح داده است. این احتمالاً به دلیل اندازه کوچک واژگان (۵۰۰ کلمه) است که باعث می‌شود رخدادهای دوکلمه‌ای پراکنده و از نظر آماری کم‌همیت باشند.

فاز ۲: ارزیابی نهایی (مجموعه آزمایشی)

با استفاده از پارامترهای بهینه یافته شده در فاز ۱ (k1=1.5, b=0.75, mu=500, lambda=0.9)، مدل‌ها را بر روی مجموعه داده آزمایشی (دیده نشده) ارزیابی کردیم.

```

58 print("\n--- Phase 2: Final Evaluation on Test Dataset ---")
59 # Constraint: Evaluate on Test using best params.
60
61 # 1. BM25 Final
62 bm25_final = BM25(VOCAB, k1=best_bm25_params[0], b=best_bm25_params[1])
63 p5, mrr, map_ = evaluate_model(bm25_final, test_questions, test_judg, test_passages)
64 print(f"BM25 TEST RESULTS -> P@5: {p5:.4f}, MRR: {mrr:.4f}, MAP: {map_:.4f}")
65
66 # 2. Unigram Final
67 uni_final = UnigramLM(VOCAB, TRAIN_COUNTS, mu=best_mu)
68 p5, mrr, map_ = evaluate_model(uni_final, test_questions, test_judg, test_passages)
69 print(f"Unigram TEST RESULTS -> P@5: {p5:.4f}, MRR: {mrr:.4f}, MAP: {map_:.4f}")
70
71 # 3. Bigram Final
72 base_uni_final = UnigramLM(VOCAB, TRAIN_COUNTS, mu=best_mu)
73 bi_final = BigramLM(VOCAB, TRAIN_COUNTS, base_uni_final, lam=best_lam)
74 p5, mrr, map_ = evaluate_model(bi_final, test_questions, test_judg, test_passages)
75 print(f"Bigram TEST RESULTS -> P@5: {p5:.4f}, MRR: {mrr:.4f}, MAP: {map_:.4f}")

```

شكل ١٦

```

--- Phase 2: Final Evaluation on Test Dataset ---
BM25 TEST RESULTS -> P@5: 0.2800, MRR: 0.4648, MAP: 0.2198
Unigram TEST RESULTS -> P@5: 0.2667, MRR: 0.3744, MAP: 0.1927
Bigram TEST RESULTS -> P@5: 0.2933, MRR: 0.4295, MAP: 0.2156

```

شكل ١٧

۶. بحث و تحلیل

مقایسه مدل‌ها

نتایج یک سلسله مراتب واضح در عملکرد را نشان می‌دهند:

۱. **BM25** به طور کلی قوی‌ترین مدل بود و بالاترین MRR و MAP را کسب کرد. توانایی آن در اشباع فرکانس کلمه (جلوگیری از اینکه سندی با ۱۰۰ تکرار یک کلمه، ۱۰۰ برابر بالاتر از سندی با ۱ تکرار امتیاز بگیرد) به آن برتری داد.
۲. **Bigram LM** در معیار MAP بسیار شبیه به BM25 عمل کرد (۰.۲۱۹۸ در مقابل ۰.۲۱۵۶). در معیار Precision@5 حتی از BM25 بهتر بود (۰.۲۹۳۳). این نشان می‌دهد که برای رتبه‌های بسیار بالا (۵ نتیجه اول)، در نظر گرفتن بافت محلی کلمات به یافتن اسناد مرتبط کمک کرده است، حتی اگر رتبه‌بندی کلی (MAP) کمی پایین‌تر باشد.
۳. **Unigram LM** از هر دو مدل دیگر عقب ماند. بدون نرمال‌سازی طول (مانند BM25) یا آگاهی از بافت (مانند Bigram)، این مدل در تمایز بین تطابق‌های باکیفیت و اسنادی که صرفاً کلمات پرسش را زیاد تکرار کرده بودند، مشکل داشت.

تنگنای واژگان

محدودیت استفاده از تنها ۵۰۰ کلمه برتر از مجموعه آموزشی احتمالاً سقف عملکرد همه مدل‌ها را فشرده کرده است.

- **پراکندگی (Sparsity)**: اصطلاحات خاص پرسش (مانند اسمی خاص یا اصطلاحات فنی) اگر در ۵۰۰ کلمه برتر نبودند، احتمالاً در طول پیش‌پردازش فیلتر شده‌اند.
- **شکست Bigram**: ترجیح برای ۰.۹ lambda در مدل Bigram این موضوع را تایید می‌کند. با چنین واژگان کوچکی، دو کلمه‌ای‌های معتبر ("new york" مثلاً) ممکن است یک یا هر دو کلمه را از دست داده باشند، که مؤلفه Bigram را بی‌استفاده می‌کند. مدل "یاد گرفت" که مؤلفه Bigram را نادیده بگیرد تا دقیق را به حداقل برساند.

شکاف تعمیم (Generalization Gap)

افت عملکرد قابل توجهی بین مجموعه اعتبارسنجی (حدود ۰.۳۲ MAP) و مجموعه آزمایشی (حدود ۰.۲۲ MAP) وجود داشت. این افت ۳۰ درصدی دو احتمال را نشان می‌دهد:

۱. عدم تطابق داده‌ها: پرسش‌های تست ممکن است ذاتاً دشوارتر بوده یا بر موضوعاتی متمرکز بوده‌اند که در واژگان ۵۰۰ تایی کمتر نمایندگی شده‌اند.

۲. بیشبرازش (Overfitting): پارامترهای تنظیم شده بر روی مجموعه اعتبارسنجی (به ویژه lambda در Bigram) ممکن است مختص ویژگی‌های پرسش‌های اعتبارسنجی بوده باشند.

نتیجه گیری

در این پژوهه، سه مدل بنیادی بازیابی اطلاعات را تحت محدودیت‌های سخت پیاده‌سازی و ارزیابی کردیم. مدل **BM25** ثابت کرد که موثرترین و قابل تعمیم‌ترین رویکرد برای این مجموعه داده است، زیرا فرکانس کلمه و طول سند را به طور موثر متعادل می‌کند. مدل زبانی **Bigram**، با وجود اینکه در $P@5$ خوب بود، از پراکندگی داده‌های ناشی از محدودیت واژگان رنج برد و مجبور شد رفتاری تقریباً مشابه با مدل **Unigram** داشته باشد.