

# سوالات بخش Index Compression

## سوال ۱

چرا در فشردهسازی postings، قبل از هر چیز docIDها را به شکل gap (شکاف) تبدیل می‌کنیم؟  
در پاسخ خودت حتماً بگو:

- چه خاصیتی در لیست postings باعث می‌شود تعریف gap معنی‌دار باشد.
- چرا برای واژه‌های پُرتکار، مقدار این gapها معمولاً کوچک است.
- کوچک بودن gapها چطور به نفع روش‌های کدگذاری مثل Gamma و Variable-Byte می‌شود.

## سوال ۲

کدگذاری Variable-Byte با اندازه پکت  $b=8$  را توضیح بده.  
در پاسخ:

1. نقش «بیت ادامه» (continuation bit) چیست و چه مقادیری می‌گیرد؟
2. ۷ بیت بعدی در هر پکت چه کاری انجام می‌دهند؟
3. از چه جهت می‌گوییم این روش از UTF-8 الهام گرفته شده است، و چرا استفاده مستقیم از UTF-8 برای فشردهسازی postings نوعی over-engineering محسوب می‌شود؟

## سوال ۳

لیست مرتب docID زیر را برای یک واژه در نظر بگیر:

[56, 58, 63, 64, 67]

1. این لیست را به لیست gap تبدیل کن.
2. الگوریتم کلی کدگذاری یک عدد مثبت  $NN$  با Variable-Byte و  $b=8$  را، مرحله به مرحله (بدون محاسبه باینری دقیق)، توضیح بده.
3. توضیح بده دیکودر چگونه می‌تواند بدون داشتن هیچ اشاره‌گر اضافی، gapها را از روی دنباله بایت‌ها بازیابی کند.

## سوال ۴

دو پیاده‌سازی متفاوت از Variable-Byte را در نظر بگیر:

- طرح A: پکتها با اندازه  $b=4$  بیت (۱ بیت ادامه + ۳ بیت داده).
- طرح B: پکتها با اندازه  $b=8$  بیت (۱ بیت ادامه + ۷ بیت داده).

به پرسش‌های زیر پاسخ بده:

1. برای چه نوع توزیع gap (اکثرا کوچک یا اکثرا بزرگ بودن gapها) طرح A از نظر تعداد بیت، فشرده‌تر از طرح B خواهد بود؟ چرا؟
2. در چه حالاتی انتظار داری طرح B از نظر سرعت دیکود (CPU کمتر) بهتر از طرح A عمل کند؟
3. این تفاوت‌ها را به طور کوتاه به مبحث کلی «تجارت بین O/I و CPU» در فشردهسازی شاخص ربط بده.



# پاسخ ها:

## پاسخ سوال ۱

- لیست postings برای هر واژه به صورت صعودی مرتب از نظر docID است، بنابراین اختلاف بین دو docID متواالی همیشه غیرمنفی و معنی دار است (gap).
- در واژه های پُر تکرار، همان واژه در سند های زیادی ظاهر می شود و این سند ها از نظر شناسه به هم «نزدیک» هستند، پس اختلاف بین docID های متواالی (gapها) اغلب عدد های کوچک مثل ۱، ۲، ۳ و ... می شود.
- وقتی gapها کوچک آند، روش های طول متغیر مثل Variable-Byte و Gamma می توانند این اعداد کوچک را با بیت های خیلی کم کد کنند، در نتیجه اندازه کلی postings به طور چشمگیری کاهش می یابد.

## پاسخ سوال ۲

1. در Variable-Byte با  $b=8$ ، بیت اول هر بایت «بیت ادامه» است: اگر باشد یعنی این بایت هنوز وسط نمایش همان عدد است و بایت های دیگری نیز برای این عدد خواهند آمد؛ اگر ۱ باشد، یعنی این بایت آخرین بایت نمایش آن عدد است.
2. ۷ بیت بعدی، بخش هایی از نمایش باینری عدد را (payload) حمل می کنند؛ وقتی چند بایت پشت سر هم برای یک عدد داریم، ۷ بیت های آنها به ترتیب کنار هم قرار می گیرند تا باینری کامل عدد به دست آید.
3. این طرح از UTF-8 الهام گرفته، چون در UTF-8 هم طول کد هر کاراکتر متغیر است و بیت های ابتدایی هر بایت نقش هدر را دارند و می گویند چند بایت مربوط به همان کاراکتر است؛ اما برای postings نیاز نداریم همه الگوهای پیچیده UTF-8 (مثل ۱۱۰xxxxx، ۱۱۱۰xxxx و ...) را پیاده کنیم، و یک نسخه خیلی ساده تر با تنها یک بیت ادامه برای ما کافی است، بنابراین استفاده مستقیم از خود UTF-8 برای docID ها، نسبت به نیاز ما، over-engineering به حساب می آید.

## پاسخ سوال ۳

1. با استفاده از تعریف gap (اولی نسبت به صفر، بقیه اختلاف با قبلی):
2.  $[56, 58, 63, 64, 67] \Rightarrow [56, 2, 5, 1, 3] [56, 58, 63, 64, 67] \Rightarrow [56, 2, 5, 1, 3]$ 
  - 3. یعنی  $3 = gap_1 = 56$ ,  $gap_2 = 2$ ,  $gap_3 = 5$ ,  $gap_4 = 1$ ,  $gap_5 = 1$ .
  - 4. الگوریتم کلی کدگذاری Variable-Byte با  $b=8$  برای عدد  $NN$  :
    - عدد  $NN$  را در مبنای ۲ منویسیم.
    - رشته باینری را از راست به چپ به گروه های ۷ بیتی تقسیم می کنیم؛ اگر گروه چپ کمتر از ۷ بیت باشد، از چپ با صفر پُر می کنیم.
    - برای هر گروه به جز گروه آخر، یک بایت می سازیم که بیت اول آن . (ادامه) و بقیه ۷ بیت همان گروه است.
    - برای آخرین گروه، بایت با بیت اول ۱ (پایان) و ۷ بیت بعدی همان گروه ساخته می شود.
    - پشت سرهم گذاشتن این بایت ها، نمایش Variable-Byte آن عدد است.
5. دیکودر دنباله بایت ها را از چپ به راست می خواند؛ برای هر عدد، تا وقتی بیت اول بایت ها . است، ۷ بیت داده را به انتهای رشته جاری اضافه می کند، و وقتی به بایتی می رسد که بیت اول آن ۱ است، می فهمد نمایش این عدد تمام شده، کل بیت های جمع شده را به عدد ددهی تبدیل می کند و به عنوان یک gap خروجی می دهد، سپس همین روند را برای عدد بعدی تکرار می کند؛ این کار نیاز به اشاره گر خارجی ندارد، چون بیت ادامه در خود بایت ها تمام اطلاعات مرزبندی را فراهم می کند.

## پاسخ سوال ۴

1. وقتی تقریبا همه gapها بسیار کوچک باشند (مثلاً اغلب در بازه ۱ تا ۷)، طرح A با  $b=4$  بهتر فشرده می‌کند، چون با ۳ بیت داده می‌تواند این اعداد کوچک را در یک پکت کوتاه نمایش دهد، و سربار کلی (یک بیت کنترل + سه بیت داده) نسبت به استفاده از ۷ بیت داده در طرح B کمتر خواهد بود.
2. وقتی بسیاری از gapها بزرگ باشند (مثلاً دهها، صدها یا هزارها)، طرح B با  $b=8$  معمولاً سریع‌تر دیکود می‌شود، چون برای نمایش یک عدد بزرگ به پکت‌های کمتری نیاز دارد؛ هر عدد بزرگ در A باید به گروه‌های ۳ بیتی شکسته شود که منجر به تعداد زیادی بایت/نیم‌بایت و کار بیشتر CPU در دیکود می‌شود.
3. طرحی مثل A بیشتر روی صرفه‌جویی حداکثری در بیت‌ها برای اعداد کوچک تمرکز می‌کند و در عوض کار CPU برای دیکود را افزایش می‌دهد؛ طرحی مثل B کمی از فشرده‌سازی بهینه فاصله می‌گیرد اما دیکود را ساده‌تر و سریع‌تر می‌کند. این دقیقاً همان trade-off کلاسیک در فشرده‌سازی شاخص است: اگر O/I و فضای ذخیره‌سازی گلوگاه باشد می‌توان بیشتر فشرده کرد و CPU را درگیرتر کرد، ولی اگر CPU محدود باشد، معمولاً طرح‌های ساده‌تر و کم‌هزینه‌تر مثل  $b=8$  اولویت دارند.