

هفت: درخت‌ها

ساختمان داده‌ها و الگوریتم

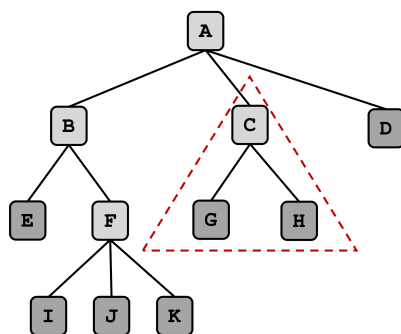
مدرس: دکتر نجمه منصوری

نگارنده: سجاد هاشمیان

۱. درخت چیست؟

درخت (tree) یک ساختمان داده‌پرا استفاده است که شبیه به یک ساختار درختی با مجموعه‌ای از گره‌های متصل به هم است؛ درحقیقت درخت یک گراف همبند بدون دور است. اکثر نویسندگان این قید را نیز اضافه می‌کنند که گراف باید بدون جهت باشد، به علاوه بعضی قید بدون وزن بودن یاها را نیز اضافه می‌کنند؛ اما لزومی به وجود این دو شرط برای تعریف درخت‌ها نیست.

۱.۱. تعاریف و انواع گره در درخت



هر گره در درخت تعدادی (صفر یا بیشتر) گره **فرزند** دارد، که در زیر آن در درخت قرار دارند (به طور قراردادی، درخت به سمت پایین رشد می‌کند، برخلاف آنچه در طبیعت می‌بینیم).

یک گره که فرزند دارد گره **پدر** آن فرزند گفته می‌شود؛ یک گره حداکثر ۱ پدر دارد. **ارتفاع** یک گره طول طولانی‌ترین مسیر پایین رو از آن گره به یک برگ است. مسیری که از گره به ریشه وصل می‌شود **مسیر ریشه** نام دارد و طول این مسیر **عمق** آن گره است.

گره‌های ریشه

بالاترین گره درخت گره ریشه نام دارد؛ پس گره ریشه پدر ندارد. این گره‌ها عملیات روی درخت معمولاً از آن شروع می‌شود. (هر چند بعضی الگوریتم‌ها از برگ شروع شده و به ریشه ختم می‌شوند). بقیه گره‌ها با دنبال کردن یاها از گره ریشه قابل دسترسی اند. در بعضی درخت‌ها، گره ریشه ویژگی‌های خاصی دارند. هر گره در یک درخت را می‌توان ریشه یک **زیر درخت** در نظر گرفت.

گره‌های برگ

پایین‌ترین گره‌های یک درخت گره‌های برگ نام دارند؛ چون این گره‌ها زیرترین گره هستند هیچ فرزندی ندارند.

گره‌های داخلی

یک گره داخلی، هر گره‌ای است که فرزند داشته باشد، پس تنها برگ‌ها گره داخلی نیستند.

۱.۲. درخت‌های منظم

به درختی گفته می‌شود که هر گره داخلی آن دقیقاً k فرزند داشته باشد؛ برخلاف تعریف عمومی درختان، در درختان k تایی ترتیب قرار گرفتن این فرزندان مهم است.

۲. نمایش درخت‌ها

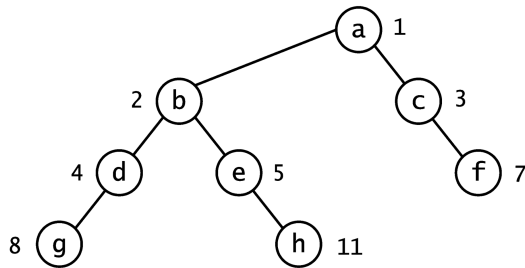
۲.۱. لیست مجاورت

همانطور که دیدیم، به طور کلی برای ذخیره سازی گراف‌ها از دو روش «ماتریس مجاورت» و «لیست مجاورت» استفاده می‌کردیم، هرچند که روش‌های دیگری نیز موجود بود، اما برای نمایش گراف‌های درختی، با توجه به این که تعداد یال‌ها از $O(n)$ است، برای استفاده از ماتریس مجاورت به یک ماتریس اسپارس می‌رسیم و این یعنی برای ذخیره کردن $O(n)$ داده باید $O(n^2)$ هزینه کنیم که اصلاً به صرفه نیست (کافیست تا n را ۱۰۰۰ در نظر بگیرید!).

۲.۲. آرایه‌ها

اما نگاهی دقیق‌تر داشته باشیم، دسته مهم و پرکاربردیی از درختان را درختان منظم تشکیل می‌دهند، که برای آنها یک فرض مازاد بر درختان داریم، یعنی برای هر راس دقیقاً k فرزند داریم که این مقدار ثابت است، پس به شکل زیر می‌توانیم عمل کنیم:

- به گره ریشه شماره (اندیس) ۰ را می‌دهیم.
- بقیه گره‌های درخت، برای هر گره با شماره v :
- به فرزند i ام این راس (در صورت وجود) شماره $kv+i$ را اختصاص می‌دهیم.

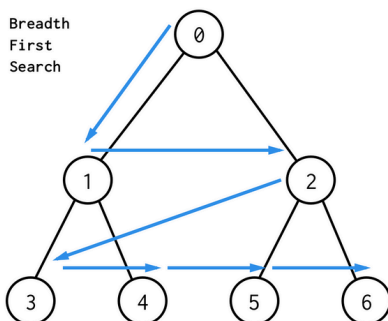


1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
a	b	c	d	e	-	f	g	-	-	h	-	-	-	-

۳. پیمایش درخت‌ها

۳.۱. پیمایش سطح اول

دوباره، البته که درخت‌ها هم یک گراف هستند و روش‌های بیان شده در باره آن‌ها هم صادق است؛ پیمایش اول سطح در درخت‌ها همانند گراف‌ها انجام می‌شود و هیچ تفاوتی ندارد، چیز متفاوت یا بیشتری را هم نشان نمی‌دهد!



با توجه به این که مرتبه زمانی این الگوریتم برابر با $O(|V| + |E|)$ است و در درخت‌ها رابطه $|E| = O(|V|)$ بنابراین این پیمایش در درخت‌ها مرتبه زمانی برابر با $O(|V|)$ دارد.

۳.۲. پیمایش عمق اول

پیمایش عمق اول نیز در درخت ها همانند گراف ها انجام می شود، اما ترتیب پیمایش راس های فرزند و ریشه برای هر زیر درخت، نکته ای است که در پیاده سازی و طراحی الگوریتم ها مورد توجه است. به طور کلی برای درختان دودویی روش های زیر توسعه داده شده اند، که به راحتی (همانگونه که در پیاده سازی ها می بینید) این روش ها به سادگی قابل تعمیم به تمامی درختان هستند.

پیمایش پیش ترتیب:

```
void preOrder(int x){
    queue.push_back(x);
    for(int i=0; i< children[x].size(); i++){
        preOrder(children[x][i]);
    }
    return;
}
```

۱. ریشه را ملاقات کن.

۲. زیر درخت چپ را پیمایش کن.

۳. زیر درخت راست را پیمایش کن.

پیمایش میان ترتیب:

```
void inOrder(int x){
    inOrder(leftChild[x]);
    queue.push_back(x);
    inOrder(rightChild[x]);
    return;
}
```

۱. زیردرخت چپ را پیمایش کن.

۲. ریشه را ملاقات کن.

۳. زیردرخت راست را پیمایش کن.

پیمایش پس ترتیب:

```
void postOrder(int x){
    for(int i=0; i< children[x].size(); i++){
        postOrder(children[x][i]);
    }
    queue.push_back(x);
}
```

۱. زیر درخت چپ را پیمایش کن.

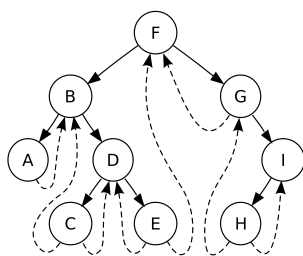
۲. زیر درخت راست را پیمایش کن.

۳. ریشه را ملاقات کن.

۳.۳. درخت نخ کشی شده (پیمایش موریس)

در یک درخت دودویی از کل اتصالات یعنی $2n$ ، تعداد $n+1$ اتصال تهی است. از این اتصالات تهی می توان برای ارتباط با دیگر گره های درخت استفاده کرد. در درختی که اتصالات تهی آن به این صورت استفاده شده است، **درخت دودویی نخ کشی شده** نامیده می شود. یک درخت دودویی را می توان به چند روش نخ کشی کرد. این نخ کشی بسته به روش پیمایش درخت دارد. برای مثال، درخت دودویی نخ کشی شده به روش میانوندی به شکل زیر تعریف می شود:

«برای هر گره، اتصال تهی سمت راست، به گره بعدی در پیمایش میانوندی اشاره کنند و اتصال تهی سمت چپ، به گره قبلی در پیمایش میانوندی اشاره کنند.»

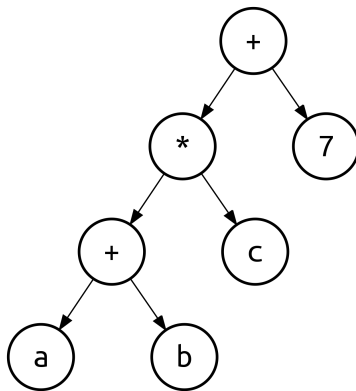


یک درخت دودویی نخ کشی شده.

درخت دودویی نخ کشی شده این امکان را فراهم می سازد که پیمایش مقادیر در درخت دودویی به روش پیمایش خطی، سریعتر از پیمایش بازگشتی درخت به صورت میانوندی صورت گیرد. همچنین پیدا کردن والد یک گره در یک درخت دودویی نخ، بدون استفاده از اشاره گرهای والد یا بدون استفاده از پشته هم امکان پذیر است، هر چند که این کار کمی آهسته است. این قابلیت وقتی مفید است که فضای پشته اندک باشد. (برای پیدا کردن گره والد به روش الگوریتم جستجوی عمق اول)

۴. کاربردهای درخت

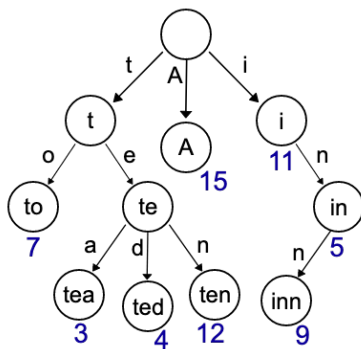
۴.۱. درخت عبارت



برای تجزیه، محاسبه و نمایش یک عبارت ریاضی (جبری، بولی، جبر مجموعه‌ها و ...) شامل ثابت‌ها، متغیرها و عملگرهای یک عملونده (قرینه، قدر مطلق، نقیض و ...) و دو عملونده (اشتراک، تقسیم و ...) می‌توان از یک درخت دودویی بهره برد که به آن درخت عبارت می‌گویند. البته برای عبارت‌هایی که عملگرهایی با بیش از دو عملوند دارند هم می‌توان درخت عبارت تشکیل داد؛ اما به ندرت چنین عملوندهایی استفاده می‌شوند و درخت حاصل دیگر دودویی نخواهد بود. برگ‌های درخت، بیانگر عملوندها (ثابت‌ها و متغیرها) هستند و هر رأس غیر برگ، یک عملگر را نشان می‌دهد.

درخت عبارت را می‌توان به سه شکل پیش‌ترتیب، میان‌ترتیب و پس‌ترتیب نمایش پیشوندی، میانوندی و پسوندی یک عبارت را تولید می‌کنند.

۴.۲. درخت پیشوندی



ترای یا درخت پیشوندی یک داده‌ساختار درختی است که برای نگاشت‌ها استفاده می‌شود (نگاشت مجموعه‌ای از زوج مرتب‌های (کلید، مقدار) است که هر کلید حداکثر یک بار می‌آید). وقتی از ترای استفاده می‌کنیم، کلیدها معمولاً رشته هستند. البته اگر نگاشتی در کار نباشد و صرفاً تعدادی رشته داشته باشیم، باز هم می‌توان برای پردازش آن‌ها از ترای استفاده کرد. گرهی ریشه نیز یک رشته‌ی خالی است. معمولاً همه گره‌ها مشخص‌کننده‌ی کلیدها نیستند. فقط برگ‌ها و بعضی از گره‌های داخلی با کلیدها مرتبط می‌شوند.

البته یک ترای الزاماً شامل رشته‌های کاراکتری نمی‌باشد؛ بلکه حتی برای جایگشت‌های عددی و مواردی از این قبیل هم می‌توان از ترای استفاده کرد.

۴.۳. کدگذاری هافمن

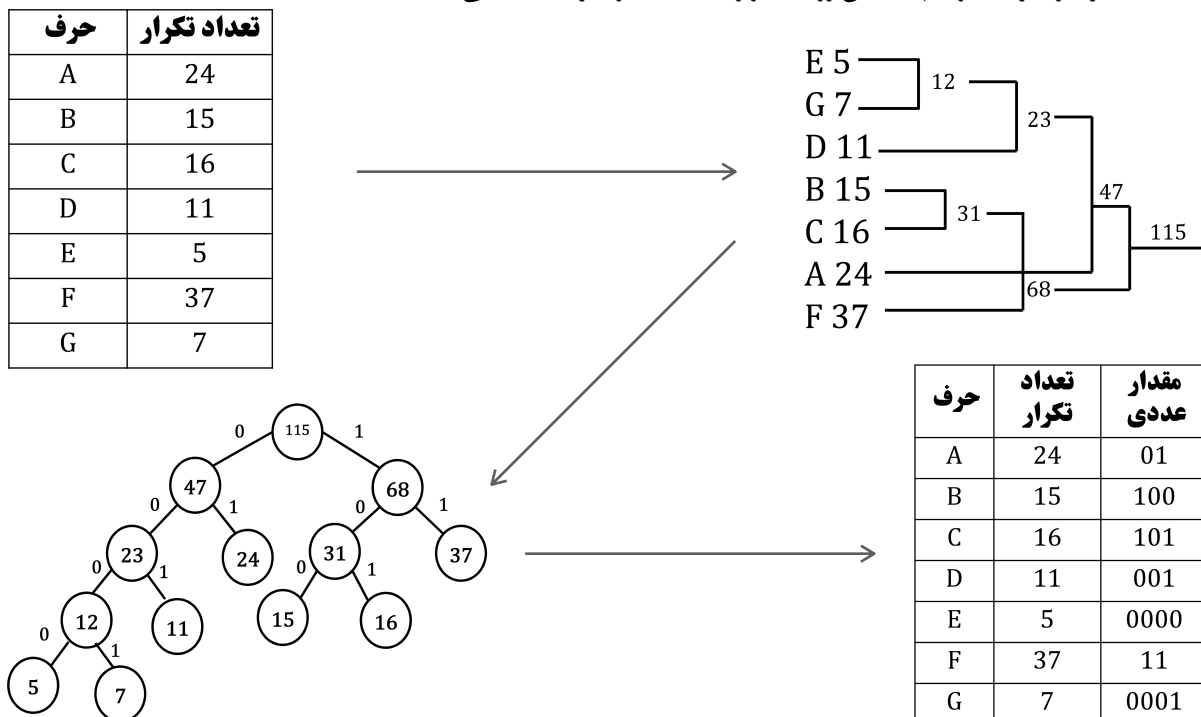
کدهای پیشوندی نوعی از کدها (توالی بیت‌ها) هستند که در آن‌ها کد اختصاص داده شده به یک کاراکتر پیشوند کد تخصیص داده شده به هیچ کاراکتر دیگری نیست. این، روشی است که کدگذاری هافمن با استفاده از آن اطمینان حاصل می‌کند که هیچ ابهامی هنگام رمزگشایی توالی بیت‌های (جریان بیت) تولید شده وجود نخواهد داشت.

فرض می‌شود که چهار کاراکتر a, b, c و d موجود هستند و کدهای طول متغیر متناظر با آن‌ها به ترتیب ۰۰۱، ۰۱ و ۱ است. این کدگذاری موجب ابهام می‌شود زیرا کد تخصیص یافته به c، پیشوند کدهای تخصیص یافته به a و b است. اگر جریان رشته فشرده شده ۰۰۰۱ است، خروجی که از حالت فشرده خارج شود امکان دارد cccb یا acd یا ab باشد.

مراحل کدگذاری برای این الگوریتم به شرح زیر است:

- I. چگالی هر کاراکتر را محاسبه میکنیم (تعداد دفعات حضور کاراکتر در متن مورد نظر).
- II. دو کاراکتر با کمترین میزان تکرار (چگالی) را انتخاب میکنیم.
- III. کاراکترهای مرحله ۲ را با کاراکتر جدیدی که دارای چگالی برابر با مجموع چگالی دو کاراکتر فوق است جایگزین میکنیم.
- IV. زمانی که فقط یک کاراکتر باقی مانده باشد، به مرحله ۲ میرویم.
- V. از عملیات فوق یک درخت حاصل می شود، بر روی این درخت هر مسیر به سمت چپ با 0 و هر مسیر به سمت راست با 1 وزن دهی میشود.

VI. کد هر کاراکتر با کنار هم گذاشتن وزن ها از ریشه تا آن کاراکتر به دست می آید.



نکته قابل توجه برای کدگذاری هافمن این است که سعی می کند با بهینه کردن میانگین طول رشته های کد شده، کمینه طول متن را برای متن کد شده داشته باشد؛ به طور دقیقتر در این الگوریتم مجموعه ای از نمادها (حروف و کاراکترها) به همراه وزن هر کدام را به ما ورودی می دهند و یک کد دودویی بدون پیشوند با کمترین امید ریاضی برای طول کد را محاسبه کند.

باید دقت کرد با توجه به کارایی بالای این الگوریتم، کدگذاری هافمن همواره بهینه نیست و در مواردی که قصد فشرده سازی بهینه تری داشته باشیم، می توان از الگوریتم های کدگذاری حسابی و یا سیستم های عددی نامتقارن استفاده کرد.

انواع مختلفی از کدگذاری هافمن وجود دارد، که بعضی از آنها از الگوریتم هایی شبیه الگوریتم هافمن و بعضی دیگر از کدهای بهینه پیشوندی (با محدودیت های خاص برای خروجی) استفاده می کنند. برای مثال می توان از کد هافمن n تایی، کد هافمن انطباقی، کد هافمن با طول محدود، کد هافمن با ارزش حرفی متفاوت و ... نام برد.

برای مطالعه بیشتر می توانید به اینجا و برای مطالعه صحت درستی این الگوریتم می توانید به اینجا مراجعه کنید.

۵. سوالات برنامه نویسی

۱. دو درخت، کوئرا

۲. تایپ خسته، کوئرا

۳. ارتباطات فامیلی، کوئرا

۴. چارلی در کازابلانکا، کوئرا

۶. برای جستجو

1. Barnes-Hut simulation
2. JSON and YAML
3. Abstract syntax tree
4. Quad Tree and Octree
5. Aho-Corasick algorithm
6. Alpha-beta pruning
7. File system
8. Lowest Common Ancestor