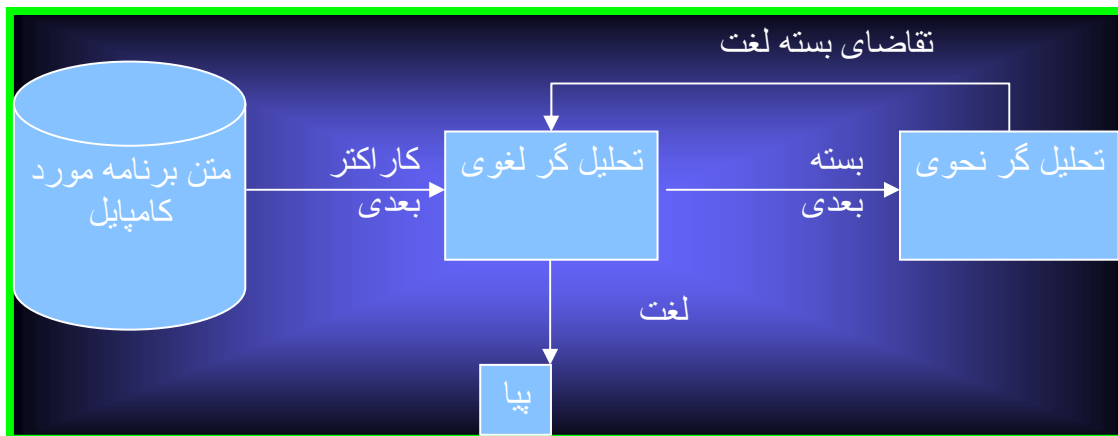


فصل دوم

تحلیلگر لغوی

۲.۱ مقدمه

تحلیلگر لغوی تابعی است که در ورودی خود برنامه مورد کامپایل را بعنوان يك فایل متن بازگشوده و کاراکتر به کاراکتر می خواند وبا رسیدن به يك کاراکتر جداکننده لغت را تشخیص و تفکیک مینماید. جداکننده ها مثل : Tab , Blank یا New line فقط به منظور جداسازی لغات بکار میروند. دسته دیگر از جدا کننده ها مثل Semicolon یا + ارزش لغوی دارند. تحلیلگر لغوی در خروجی خود اطلاعات مربوط به لغت را در يك رکورد یا ساختار بنام بسته لغت قرار می دهد بسته لغت شامل اطلاعاتی در مورد لغت تشخیص داده شده توسط تحلیلگر لغوی مانند سطر ، ستون و نوع لغت است بسته لغت را در اصطلاح Token گویند.



شکل 2.1- شمایی کلی تحلیلگر لغوی

هر زبان برنامه سازی قوانین لغوی خاص خود را دارد. قوانین لغوی زبانها را میتوان بصورت غیر مبهم و خلاصه در قالب نوعی خاص از عبارات بنام عبارت با قاعده بیان نمود. قوانین لغوی را نیز میتوان توسط نوعی خاص از ساختارگراف بنام ماشین های خودکار مطرح کرد. ماشینهای خودکار ابزاری برای تبدیل قوانین لغوی به کد برنامه هستند. چنانچه لغتی بنا بر قوانین لغوی زبان برنامه سازی ایجاد نشده باشد ، تحلیلگر لغوی اعلام خطا مینماید.

۲.۲ ساختار ورودی / خروجی

در داخل تابع تحلیلگر لغوي برنامه مورد کامپایل به عنوان يك فایل متن دسترسي ميشود. اين فایل در ابتدای برنامه کامپایلر و خارج از محیط تحلیلگر لغوي تعريف و گشوده مي شود و به عنوان يك پارامتر به تابع تحلیلگر لغوي ارسال مي گردد. در قطعه کد ارائه شده در شکل 2.2 چگونگی فراخوانی تابع تحلیلگر لغوي نشان داده شده است..

```

Void main ( int arg c , char *argv [ ] )
{
File *InText ;
if ( argc < 2 )
{
clrscr ( ) ;
textcolor ( RED ) ;
gotoxy ( 10 , 5 )
printf ( " نام برنامه مورد فراموش شده است " ) ;
if ( !getch ( ) ) getch ( ) ;
}
InText = fopen ( argv [ 1 ] , " R " ) ;
while ! feof ( InText ) ;
{
Token-type Token ;
Token = Lexer ( InText ) ;
}
}

```

شکل 2.2- فراخوانی تابع تحلیلگر لغوي

نقطه گنگ در مثال فوق TokenType و در واقع نوع خروجی تحلیلگر لغوي است. معمولاً اطلاعات زیر در بسته Token قرار داده مي شود :



```

Typedef struct Token
{ int ROW ; // شماره سطر
  int COL; // شماره ستون
  int BLKNO ; // شماره آشیانه اي
  int BLKORD ; // شماره ترتيب بلاك در برگیرنده
  enum Symbols type ; // نوع لغت
  char Name [ 30 ] ; // خود لغت
} Token-type ;

```

شکل 2.7- ساختار بسته لغت

در ساختار ارائه شده برای تعریف لغت که در اینجا بسته لغت یا Token نامیده شده است شماره سطر و ستون ظاهر شده است. با استفاده از شماره سطر و ستون لغت ، در هنگام صدور پیام خطا کامپایلر قادر به بیان دقیق مکان خطا خواهد بود. BLKNO شماره آشیانه اي بلاك در بر گیرنده يك لغت را مشخص مي كند. تنها نام لغات برای تعیین و تشخیص آنها از یکدیگر کافی نیست بلکه ، جهت تعیین يك لغت نیاز به شماره آشیانه اي بلاك در برگیرنده آن نیز هست . برای نمونه به قطعه کد زیر توجه کنید.

```

{
  int I ;
  I=5 ;
  { int I ;
    I=6 ;
    printf ( " 2nd blk %d " / I ) ;
  }
  printf ( "\n 1st blk %d " / I ) ;
}

```

در قطعه کد فوق همانگونه که مشاهده می کنید اگرچه که برای دو متغیر يك نام I تعیین شده اما شماره آشیانه اي بلاك در بر گیرنده ، وجه تمایز این دو می باشد. گاهی اوقات شماره آشیانه اي بلاك نیز کافی نیست برای نمونه در قطعه کد زیر اگرچه که شماره آشیانه اي بلاكها یکی است اما دو متغیر وجود دارد.

```

{ int I ;
  I = I + 1
}
{ int I ;
  I = I + 3
}

```



شماره آشیانه اي با ورود به بلاك افزایش و با خروج از بلاك کاهش مي يابد. با ظهور هر بلاك جديد شماره بلاك يا BLKORD يك واحد اضافه مي شود در مثال فوق وجه تمايز دو متغیر شماره ترتیب بلاك در بر گیرنده آن است .

به هر لغت يك نوع تخصیص داده میشود. نوعي شمارش پذیر يا در اصطلاح Enumerated. براي نمونه تحلیلگر لغوي اگر در متن فایل ورودی عدد 123 را تشکیل دهد، در داخل فیلد Name عدد 123 و در داخل فیلد Type، نوع آن که يك عدد صحیح مي باشد را مشخص مي کند. اگر در متن ورودی، جمله زیر قرار گرفته باشد:

ABC := 123 ;

تحلیلگر لغوي با اجرای دستور العملي مثل :

NextChar =getc (in-text);

ابتدا کاراکتر A و سپس B و بعد از آن C را از متن فایل ورودی خوانده، حالا با مشاهده Blank و در صورت عدم وجود Blank با مشاهده علامت کولون (:)، خاتمه اولین لغت را تشخیص داده، در داخل فیلد Name رشته ABC و در داخل فیلد Type، نوع آن که يك شناسه مي باشد را قرار مي دهد.

در فراخوانی بعدی تحلیلگر لغوي با Blank مواجه مي شود. از آن چشم پوشی مي کند و به جلو مي رود و علامت := را به عنوان لغت بعدی تشخیص داده در فیلد Name، رشته := و در فیلد Type، عملگر تخصیص را قرار مي دهد. البته شماره سطر وستون و بلاك در بر گیرنده هر لغت نیز مشخص مي شود.

نوع لغات در فیلد Type از ساختار TokenType مشخص شده است. این فیلد از نوع شمارش پذیر بنام Symbols تعریف شده است. در داخل نوع شمارش پذیر يا در اصطلاح Enumerated بنام Symbols انواع ممکن لغات در داخل زبان مورد نظر مشخص میشود. باید توجه داشته باشید که با تعریف متغیر از نوع شمارش پذیر مجموعه اي از مقادیر ثابت يا Constant تعریف میشود. نوع Symbols در زبان C بصورت زیر برای نمونه تعریف میشود.

```
enum Symbols { S_Program, S_Const, S_Eq, S_Semi, S_Id, S_No, S_Type,
              S_Record, S_End, S_Int, S_Real, S_Char, S_Array, S_String,
              S_Begin, S_Var, S_Colon, S_ParBaz, S_ParBast, S_Set, S_of,
              S_BrackBaz, S_BrackBast, S_Case, S_Pointer, S_ConstString,
              S_Function, S_Procedure, S_Begin, S_Dot, S_Comma,
              S_If, S_Then, S_Else, S_While, S_Do, S_Repeat, S_Until,
              S_For, S_Add, S_Sub, S_Div, S_Mul, S_Mod, S_Lt, S_Le,
              S_Gt, S_Ge, S_Gt, S_Ne, S_Not, S_And, S_Or};
```

همانگونه که مشاهده میشود انواع مختلف لغات باید در نوع Symbols گنجانده شود.

۲.۳ عبارات با قاعده

عبارت با قاعده¹ ، فرمی است چکیده و خلاصه برای بیان قوانین لغوي زبانها. شکل لغات در زبانهای برنامه سازی دارای فرم کلی خاص است. بر اساس این فرم کلی است که انواع لغات از قبیل شناسه ها (اسامي) ، اعداد ، رشته های ثابت کاراکتری ، جملات تفسیری (Comment) و سایر لغات از یکدیگر تفکیک و تمیز داده می شوند.

۲.۳.۱ نمونه هایی از عبارات با قاعده

به عنوان نمونه تعریف شناسه ها در زبان C را در نظر بگیرید. يك شناسه یا Identifire در زبان C حتماً باید با يك کاراکتر آغاز گردد و در ادامه ممکن است به هر تعداد وبا هر ترکیبی از ارقام صفر تا نه (0-9) ، حروف الفبای انگلیسی (A-Z) و علامت خط زیر یا در اصطلاح (UnderLine) ادامه یابد. به عنوان مثال اسامي :

A , B__1 , BAC2345__

از لحاظ زبان C به عنوان اسامي (شناسه) ، شناخته می شوند. میتوان با استفاده از يك عبارت با قاعده به صورت زیر ، فرم کلی شناسه ها را در زبان C تعریف نمود :

Identifire : Letter (Letter | Digit | ' _ ') *

Letter : A | B | | Z | a | b | | z

Digit : 0 | 1 | 2 | | 9

در عبارت ارائه شده برای شناسه ها ، علامت '*' ، نمایانگر تکرار صفر یا بیشتر می باشد و علامت ' | ' ، علامت یا می باشد. می توان بصورت چکیده نیز Digit را تعریف نمود :

Digit : [0 ... 9]

به همین ترتیب ، حروف الفبای انگلیسی را بصورت زیر می توان تعریف نمود :

Letter : [a ... z , A ... Z]

شکل کلی اعداد صحیح را میتوان با استفاده از يك عبارت به این صورت مشخص نمود:

Number : digit digit *

Digit : [0 ... 9]

بر طبق این عبارت ، يك عدد صحیح با يك رقم بین صفر تا نه آغاز می شود و به هر تعدادی رقم می تواند در ادامه آن ظاهر شود. برای مثال :

0 , 5 , 0123



همانگونه که مشخص است ، يك عدد حداقل داراي يك رقم بايد باشد يا به عبارت ديگر يك عدد از يك يا بيشتر ارقام تشكيل شده ، مي توان عدد را بصورت زير نيز تعريف کرد :

$$\text{Number} : \text{digit} +$$

در اينجا علامت '+' نمايانگر تکرار يك يا بيشتر مي باشد، يادآوري مي كنيم كه علامت ستاره '*' نمايانگر تکرار صفر يا بيشتر است.

اعداد مي توانند داراي علامت و يا بدون علامت باشند

$$\text{Number} : (+ | - | \epsilon) \text{digit} +$$

در اين عبارت با قاعده ، علامت اپسيلون ϵ نمايانگر عدم وجود و ياتهي مي باشد . عبارت فوق به اين صورت خوانده مي شود :

يك عدد داراي علامت بعلاوه يا منها و يا ممكن است اصلاً "علامتي نداشته باشد و بدنبال آن به تعداد يك يا بيشتر ارقام بين صفر تا نه ظاهر مي گردد.

اعداد را ميتوان به فرمي ساده تر با استفاده از اين قاعده كه هر عبارت اختياري مثل r بصورت $r?$ نمايش داده ميشود ، بصورت زير مشخص نمود :

$$\text{Number} : (+ | -) ? \text{digit} +$$

بعنوان نمونه اي ديگر از عبارات با قاعده ، فرم كلي رشته ها در زبان پاسكال را در نظر بگيريد. بايد توجه داشته باشيد كه در اين زبان چنانچه در سطح رشته علامت کوتیشن ' نياز باشد ، مي بايست آنرا دوبار تکرار نمود. براي نمونه چنانچه جمله 'This is your's در خروجي مورد نظر باشد، ميتوان دستور ' Writeln (' s ') This is your استفاده نمود. لذا، فرم كلي رشته ها با عبارت با يك عبارت با قاعده بصورت زير بيان ميشود :

$$\text{String} : ' (\text{Characters} | ') * ' '$$

در عبارت فوق ، مجموعه Characters نمايانگر هرگونه کاراکتر قابل مشاهده منهاي کوتیشن است.



برای ایجاد یک عبارت با قاعده، تعدادی از علائم مورد استفاده واقع می شوند. هر یک از این علائم دارای معنی و مفهوم خاصی می باشد. در حالت کلی اگر دلتا (Δ) مجموعه علائم بکار رفته شده در عبارات باشد، آنگاه:

1. هر عنصر $a \in \Delta$ خود یک عبارت با قاعده است. برای مثال چنانچه $\Delta = [0..9]$ باشد آنگاه هر رقمی بین صفر تا نه مثل 1، خود به تنهایی یک عبارت با قاعده است.

2. چنانچه r و s دو عبارت با قاعده باشند، آنگاه rs نیز یک عبارت با قاعده است و به عبارت ساده تر در حالت کلی اگر دو عبارت با قاعده را در کنار یکدیگر قرار دهید، حاصل یک عبارت با قاعده خواهد بود.

3. چنانچه r و s دو عبارت با قاعده باشند، آنگاه r یا s که بصورت $(r|s)$ مشخص می شود نیز یک عبارت با قاعده است و عملگر $|$ دارای خاصیت جابجایی است، به عبارت دیگر:

$$r|s = s|r = (r|s)$$

4. چنانچه r یک عبارت با قاعده باشد آنگاه r^* نمایانگر تکرار صفر یا بیشتر از عبارات r است و برای نمونه اگر r ، a یا b باشد آنگاه r^* برابر است با:

$$r : a|b \\ : (a|b)^*$$

این عبارت گویای هر ترکیبی با هر تعدادی از a یا b که در کنار یکدیگر قرار گرفته اند می باشد. برای نمونه رشته های
AAA, BAABB, BBBB

همگی فرم های خاصی از عبارت r^* هستند بنابراین علامت اپسیلون ϵ که نمایانگر عنصر تهی است، به تنهایی یک عبارت با قاعده است

5. چنانچه r یک عبارت با قاعده باشد، آنگاه r^+ نمایانگر تکرار یک یا بیشتر عبارت r است. برای نمونه:

$$\text{Number} : \text{digit}^+$$

$$r^+ = (r^*|\epsilon)$$

با استفاده از پنج قاعده فوق، می توان عبارات با قاعده را بنا نمود.

با استفاده از نکات فوق، فرم کلی جملات تفسیری (Comment) در زبان پاسکال را میتوان بصورت زیر معین نمود. باید توجه داشته باشید که در زبان پاسکال جملات تفسیری در بین علائم آکولاد باز و آکولاد بسته قرار می گیرند.

$$\text{Comment1} : \{C^*\}$$

$$C : (\text{کلیه کاراکترها منهای آکولاد})$$

در زبان پاسکال جملات تفسیری را میتوان بین علائم $(^*$ و $)^*$ نیز مشخص نمود. در این فرم از جملات تفسیری باید توجه داشته باشید که اگر ستاره در داخل جمله



ظاهر شود، در پشت آن باید هر حرفي به غير از ستاره و پرانتز بسته ظاهر شود
براي نمونه ، لغت زیر يك جمله تفسيري نیست :

(* abc *) acb *)

فرم كلي اینگونه جملات را میتوان در قالب يك عبارت با قاعده بصورت زیر خلاصه نمود:

Comment2 : ‘(* ((r | *⁺s) * r) * *⁺)’

r : کلیه کاراکترها منهای علامت *

s : کلیه کاراکترها منهای علائم * و ‘

به این ترتیب جملات تفسيري بصورت زیر

Comment : Comment1 | Comment2

تعريف مي شوند.

با توجه به اینکه اعداد یا به صورت صحيح و یا به صورت اعشاري ظاهر مي شوند و با در نظر گرفتن اینکه اگر قبل از ممیز حداقل يك رقم ظاهر شود آنگاه وجود ارقام پس از ممیز ضروري نیست و نیز اینکه اگر بعد از ممیز حداقل يك رقم ظاهر شود وجود رقم قبل از ممیز ضروري نیست . بطور خلاصه عبارت با قاعده براي بيان شکل كلي اعداد بصورت زیر میتواند باشد:

Number : (digit*.digit⁺) | (digit⁺.digit*) | digit⁺

عبارات با قاعده مبین قوانین لغوي و نمایانگر فرم كلي لغات هستند . متأسفانه ، این فرم كلي را بسادگی نمیتوان تبدیل به کد برنامه نمود . لذا ، براي رفع این مشکل از ماشینهای خودکار استفاده میشود .

۲.۵ - ماشینهای خودکار

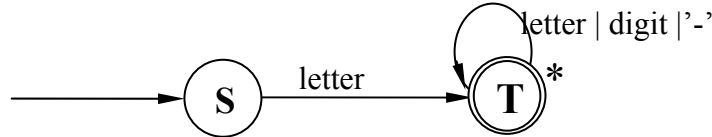
يك ماشین خودکار² نوعي گراف مي باشد که دارای تعدادي رئوس یا حالات (States) و تعدادي لبه یا یال (Edges) مي باشد و در این گراف یالها خطوط واصل بین حالات هستند . این ماشینها، ابزاري براي تعيين قوانین لغوي و تشخیص لغات مي باشند که بسادگی قابل تبدیل به کد برنامه بوده بطوري که با استفاده از آنها مي توان تحلیلگر لغوي را بصورت يك برنامه تبدیل کرد و یا Source يك برنامه را تولید نمود.

براي نمونه شناسه ها را در نظر بگیرید . شناسه ها توسط عبارت با قاعده :



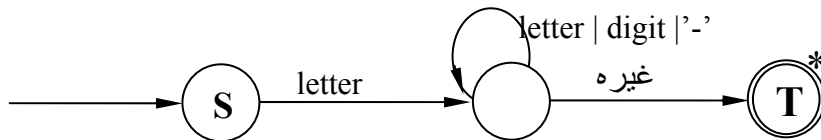
letter (letter | digit)*

در بخش قبل معین شدند. برای تشخیص شناسه ها می توان ماشین خودکار زیر را مورد استفاده قرار داد:



این ماشین خودکار را بصورت زیر نیز میتوان نمایش داد. بنابر دیاگرام فوق در صورتیکه کاراکتر خوانده شده از متن فایل مورد تحلیل لغوي یکی از حروف الفبای لاتین یا در اصطلاح letter باشد میتوان وارد ماشین خودکار تشخیص شناسه ها شد. حالت شروع و در واقع نقطه ورود به ماشین خودکار با حرف S مشخص شده است. حالت بعدی T نامیده شده است. در حالت T اطمینان از تشخیص يك شناسه است زیرا، يك حرف نیز به تنهایی يك شناسه در نظر گرفته میشود. پس در این حالت لغت خوانده شده بعنوان يك شناسه پذیرفته شده است.

حالت پذیرش ماشین خودکار را با دو دایره تو در تو نشان میدهند. در حالت پذیرش T آنقدر از ورودی کاراکتر های بعدی خوانده میشوند تا به کاراکتری غیر از آنچه بر روی این حالت مشخص شده در ورودی برسید. پس برای تشخیص يك شناسه نیز نیاز بخواندن يك کاراکتر اضافی از ورودی هست. در شکل فوق علامت * نمایانگر خواندن يك کاراکتر اضافی تر است. برای مثال در <AB12 با دیدن علامت < است که می توان فهمید شناسه AB12 به پایان رسیده است. میتوان شناسه ها را با ماشین خودکار زیر نیز نمایش داد:



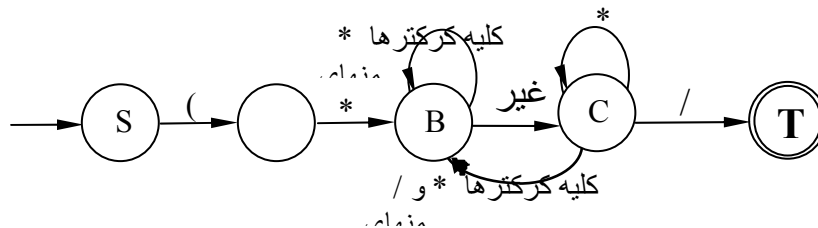
جملات تفسیری Comment توسط تحلیلگر لغوي بعنوان يك لغت شناسایی میشوند. جملات تفسیری پاسکال را میتوان توسط عبارت زیر مشخص نمود:

Comment2 : ‘(* ((r|*s)*r)* *+’

r : کلیه کاراکترها منهای علامت *

s : کلیه کاراکترها منهای علائم * و ‘

میتوان جملات تفسیری را با يك ماشین خودکار بصورت زیر مشخص نمود:



همانگونه که مشاهده می کنید در حالت شروع S با دیدن کاراکتر / ، گذری (Transision) از حالت شروع S به حالت A وجود دارد و در حالت A اگر * دیده شود گذری به حالت B وجود دارد. در غیر اینصورت این ماشین خودکار لغت خوانده شده را نمی پذیرد و در اصطلاح Reject می کند. در حالت B با دیدن هر کاراکتری غیر از * ماشین خودکار به کار خواندن کاراکترهای بعدی ادامه می دهد. اما ورودی * اهمیت دیگری دارد. با دیدن * اگر در پشت آن علامت / ظاهر شود کار خاتمه یافته تلقی می شود.

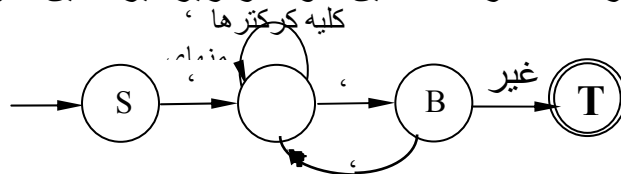
حالت T ، حالت پذیرش است. در حالت C در واقع ماشین بخاطر دارد که * دیده شده است. اگر در حالت C در ورودی علامت / ظاهر شود ، کار خاتمه یافته است. در غیر اینصورت ، به حالت B تغییر وضعیت داده میشود. پس مشاهده می کنید که ماشینهای خودکار بسیار ساده تر از عبارات با قاعده قابل تولید هستند

فرم کلی رشته ها را در زبان پاسکال با عبارت با قاعده

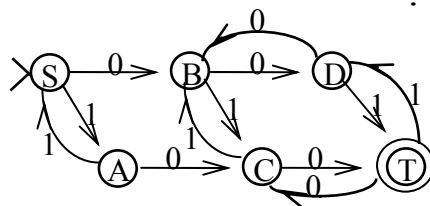
String : '(c|')*

کلیه کاراکترها منهای علامت کوتیشن ' C :

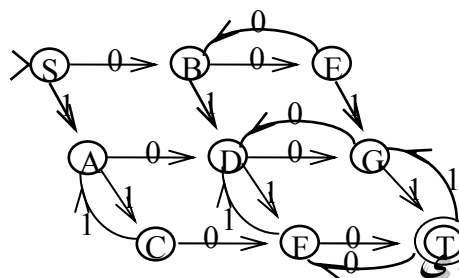
مشخص میشود. میتوان رشته ها را با ماشین خودکار زیر نیز معین نمود:



مثال: يك ماشین خودکار قطعی برای کلیه اعداد مبنای دو که تعداد صفرهای آنها زوج و تعداد یک ها فرد است ایجاد کنید. در ضمن لااقل دو عدد صفر و يك عدد يك باید در این اعداد موجود باشد.



مثال: يك ماشین خودکار قطعی برای کلیه اعداد مبنای دو که تعداد صفرهای آنها زوج و تعداد یک ها زوج است ایجاد کنید. در ضمن لااقل دو عدد صفر و دو عدد يك باید در این اعداد موجود باشد.



۲.۶ ایجاد تابع تحلیلگر لغوی

شکل 2.8 - دیاگرام تابع تحلیلگر لغوي

هرگاه که تابع تحلیلگر لغوي مورد فراخواني واقع مي گردد ، در حالت شروع صفر قرار مي گيرد . اگر در فراخواني قبل، کاراکترهاي اضافه خوانده شده بود ، آن کاراکتر اضافه را مورد استفاده قرار مي دهد ، وگرنه کاراکتر بعدي را از داخل متن برنامه مورد کامپایل مي خواند. تا زماني که کاراکترهاي که ارزش لغوي ندارند مانند Tab ، Blank و Newline خوانده شوند ، تحلیلگر لغوي در همان حالت شروع صفر باقي ميماند. در غير اينصورت وابسته به نوع کاراکتر ، در يك جمله Case اقدام به تشخيص لغات مختلف مي نمايد .

میتوان با استفاده از دیاگرام شکل 2.8 کد تابع تحلیلگر لغوي را بشرح زیر در زبان C ايجاد نمود.

```

struct TokenType lexer ( FILE *InText
)
{
    enum Symbols LexiconType ;
    char NextChar ; NextWord[80] ;
    int State , Lentgh ;
    static char LastChar = '\0';
    static int RowNo =0, ColNo = 0;
    State = 0 ; // وضعیت شروع قرار مي گيرد
    Length = 0 ;
    while ( ! feof ( InText ) )
    {
        if ( LastChar )
            { Nextchar = LastChar ; LastChar = '\0' ; }
        else
            NextChar = fgetc ( InText ) ;
            NextWord[Length++] = NextChar;
            switch State
            {
                case 0: // حالت شروع صفر
                    if (NextChar == '\n') {RowNo++; ColNo = 0; }
                    else ColNo++;
                    if ( Nextchar == ' ' || Nextchar == '\t ' ||
                        Nextchar == '\n ' ) Length = 0;
                    else if (( Nextchar <= ' z ' &&
                        Nextchar >= ' a ' ) ||
                        ( Nextchar <= ' Z ' &&
                        Nextchar >= ' A ' ) State = 1;
                    else if ( Nextchar <= ' 9 ' && Nextchar >= ' 0 ' )
                        State = 2 ;
                    else if ( Nextchar == ' ( ' ) State 3;

```



```

else if ( Nextchar == '<' )      State 4
else if ( Nextchar == '>' )      State 5
else  LexerError(NextWord, Length);
break ; // پایان حالت صفر
case 1 : // تشخیص شناسه ها
    if( Isalpha( Nextchar ) || Isdigit( Nextchar ) ||
        NextChar == '_' ) State = 1;
    else { Lastchar = Nextchar ;
          NextWord[Length-2] = '\0';
          return MakeToken(IsKeyWord(NextWord)); }

Break ;
case 3 : // تشخیص اعداد
    .
    .
    .
} // پایان سوئچ
} // پایان حلقه

```

در مثال فوق تابع MakeToken کار ایجاد بسته لغات یا در اصطلاح Token را بر عهده دارد. IsKeyWord کار تشخیص نوع لغات کلیدی را بعهدده دارد. این تابع در زبان C بشرح زیر است:

```

enum Symbols IsKeyWord( char *key)
{
    int I ;
    struct KeyType
    { char *key;
      enum Symbols Type
    } KeyTab[] = { "if", S_If,
                  "while", S_While,
                  "then", S_Then,
                  "else", S_Else,
                  "integer", S_Integer,
                  "type", S_Type,
                  "function", S_Function,
                  0, 0};
    for( I =0; KeyTab[ I ].key &&
          strcmp(KeyTab[I].key, Key); I++ ) ;
    if(KeyTab[I].key) return KeyTab[I].Type;
    Return S_Identifier;
} // EOF IsKeyWord

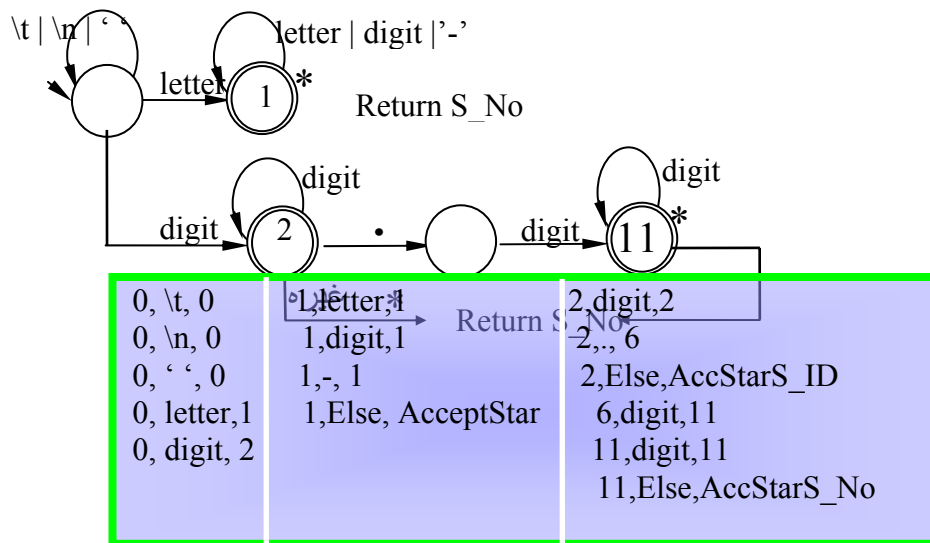
```



میتوان به سادگی تابع مولد تحلیلگر لغوي ایجاد نمود. کافیت ، تابعي براي پیمایش يك گراف بنویسید. این تابع در ورودي خود علاوه بر متن برنامه مورد تحلیل لغوي ، فایل قوانین لغوي را در فرم يك ماشین خودکاري پذیرد. این تابع با خواندن هرکاراکتر از داخل فایل متن برنامه مورد تحلیل لغوي به پیمایش ماشین خودکار میپردازد تا نوع لغت را تشخیص دهد.

از مولدهای شناخته شده تحلیلگر لغوي ، يکي LEX میباشد. این مولد امروزه بخصوص بر روي سیستمهاي عامل UNIX موجود مي باشد. کار با LEX مشکل است چرا که نیاز به بیان قوانین لغوي ، در فرم عبارات با قاعده دارد. در صورتی که همانگونه که تا کنون مشاهده کرده اید، به فرمی خیلی ساده تر می توان با استفاده از ماشینهاي خودکار ، قوانین لغوي را بیان نمود .

باید قوانین لغوي را در قالب يك ماتریس مشخص نموده و در داخل يك فایل متن قرار داد. سپس در داخل برنامه میتوان به این فایل متن ارجاع و ماشین خودکار را در قالب يك ماتریس مشخص نمود. برای نمونه در زیر قوانین لغوي در قالب ماشین خودکار و ماتریسی مشخص شده است.





برای ایجاد ماتریس نیاز به یک زبانی ساده است که بر اساس آن برای نمونه بتوان مشخص نمود که برای مثال Letter نیست و یا غیره ، که می باشد . شاید بهتر بود که در فایل ورودی در داخل ماتریس قوانین لغوي Letter بصورت زیر مشخص میشد:

0 ; [A - Z , a - z] ; 1

حالا در داخل برنامه ، این ساختار را به عنوان یک فایل TEXT ، کاراکتر به کاراکتر خوانده ، عینا در داخل یک ماتریس با تعداد سطرها که مساوی با تعداد رکوردهای فایل است و تعداد ستونها که مساوی با طول هر رکورد یا هر سطر است ، ضبط نمود.

اکنون یا با نوشتن یک برنامه پیمایش کننده ماشین خودکار و یا تولید متن یک برنامه C که بر اساس این ماشین خودکار برای انجام عمل تحلیل لغوي ایجاد شده ، می توان لغات را تشخیص داد . همزمان با خواندن هر کاراکتر از ورودی ، مثلاً با خواندن کاراکتر **کوجکتر** از ورودی طبق ماتریس از حالت شروع صفر به حالت 5 تغییر وضعیت داده می شود ، در اینجا سرعت عملیات تحلیلگر لغوي ، وابسته به سرعت جستجو در داخل ماتریس است .

می توان با ایجاد یک تحلیلگر لغوي برای تشخیص متن برنامه ها با لغات کلیدی فارسی اقدام نمود . برای مثال جمله زیر را میتوان به یک جمله if تبدیل نمود .

حقوق > 1000

اگر

: 100 / حقوق = مالیات آنگاه

۲.۸ تبدیل برنامه ها به زبان فارسی

میتوان متن برنامه های فارسی را که بر اساس قواعد زبانهایی متفاوت نوشته شده است را با استفاده از یک جدول تبدیل کلی به زبان انگلیسی تبدیل و بالعکس پس از کامپایل و اشکال زدایی برنامه ها دوباره با استفاده از همان جدول از انگلیسی به فارسی تبدیل نمود .

اصولاً زبانهایی سطح بالا به خاطر خوانا بودنشان ، **سطح بالا** خوانده می شوند . متأسفانه این انگلیسی بودن زبانها، خوانایی برنامه ها را با مشکل مواجه نموده است . بخصوص در انتهای مراحل تجزیه و تحلیل سیستم ها، برنامه سازان را در تبدیل متن عملیات تحلیل شده از شبه دستورالعمل یا در اصطلاح Pseudo code به کد برنامه عملاً با مشکل زیاد روبرو نموده است . شبه دستورالعمل را فارسی زبانها نمی توانند به سادگی به زبان انگلیسی بنویسند . شبه دستورالعمل ها معمولاً از زبانهایی سطح بالا الهام گرفته می شوند . مشکل زبان فارسی از راست به چپ بودن جملات و بالعکس ، چپ به راست بودن عبارات است . البته ، در اینجا این مشکل موردی ندارد .



در مورد شناسه ها نیز در زبان فارسي مشکل وجود دارد . شناسه ها از انتها تشخیص داده مي شوند . مشکل ديگر اين است که در زبان انگلیسی 28 و در زبان فارسي 32 حرف وجود دارد . البته ، با ترکیب حروف مي توان اين مشکل را حل نمود .

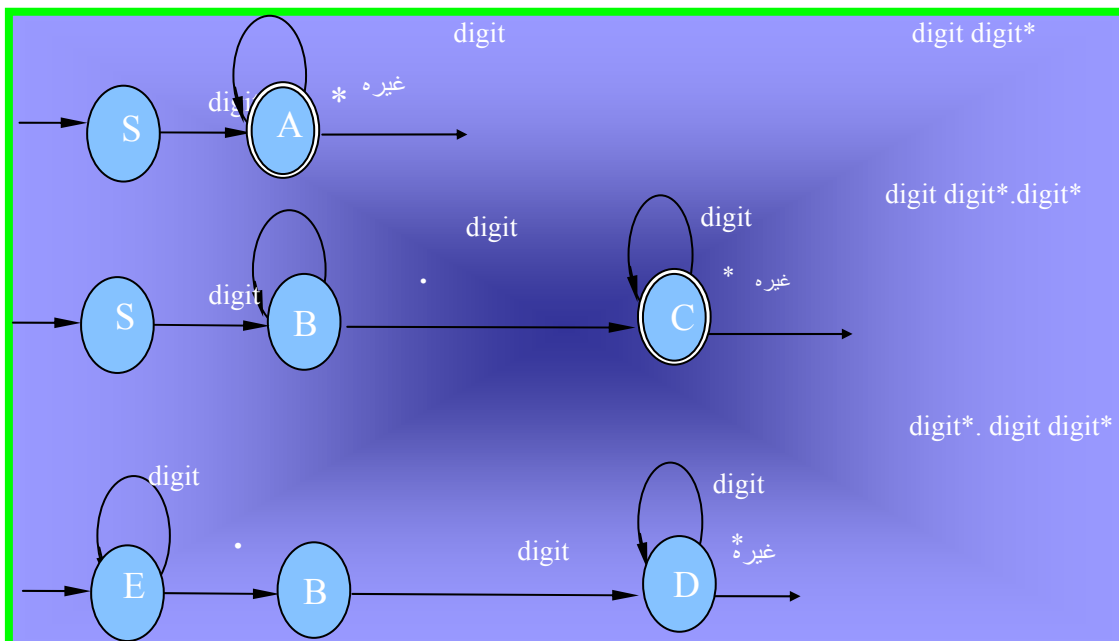
نکته ديگر ، محیط ویرایش متن برنامه هاست . در اغلب محیط هاي ارائه شده براي کامپایلرها ، **Syntax Directed Editors** ویرایشگرهايي مي باشند که بنابر قوانین زبان به برنامه ساز کمک مي کنند تا بتواند برنامه خود را با Syntax صحيح ایجاد کند . ميتوان يك محیط قوي ویرایش وابسته به زبان مورد نظر براي ویرایش برنامه هاي فارسي ایجاد نمود .

به اين ترتيب ميتوان برنامه ها را به زبان فارسي با استفاده از يك ویرایشگر باهوش براي زبان خاص نوشت . سپس با استفاده از جدول تبدیل برنامه به زبان اصلي ترجمه و پس از اشکالزدایی دوباره از زبان اصلي با استفاده از جدول تبدیل لغات را به فارسي برگرداند .

۲.۹ ماشینهای خودکار غیر قطعی

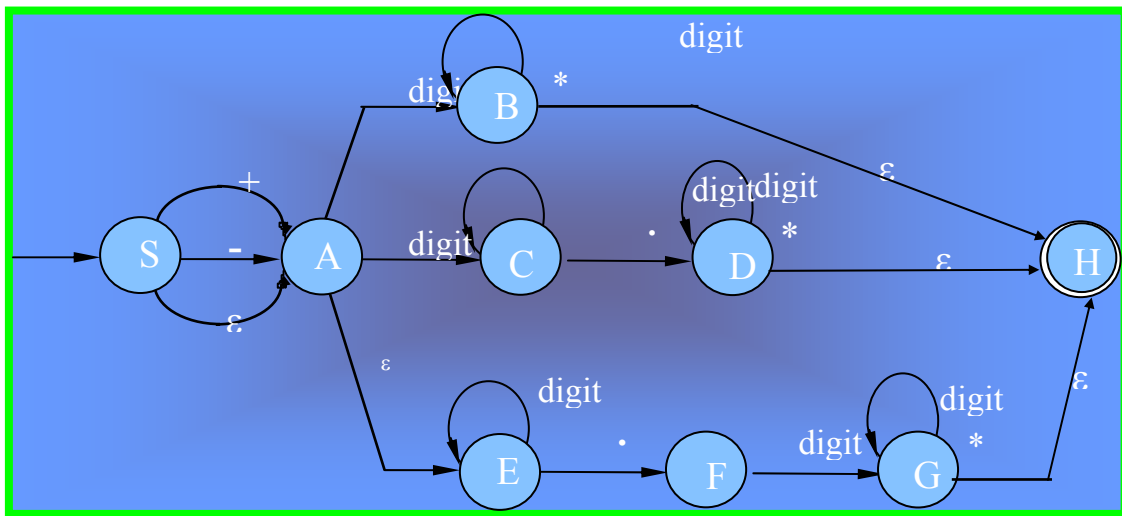
اصولا يك ماشين خودكار در واقع ابزاري براي نمايش چگونگی تصميم گيريهاست . ماشين خودكار يا بصورت قطعي تصميم گيريهي را انجام مي دهد يا غير قطعي. اگر به ماشينهاي خودكاري كه تا به حال ترسيم شده توجه كنيم ، مشخص است كه در هر حالي بطور قطعي مشخص ميكنند كه به ازاي ديدن يك رشته در ورودي از يك حالت به کدام حالت مي بايست گذر نمود . ماشينهاي خودكار قطعي را اصطلاحاً "Determistic Finite State Automata" گویند.

در ماشينهاي غير قطعي يا در اصطلاح NonDeterministic ، عدم قطعيت در تصميمات وجود دارد . البته اين عدم قطعيت اگرچه كه كار تصميم گيري را مشكل ميکند، اما كار ايجاد ماشين را بسيار ساده تر مي نمايد. حالا، نکته در اينجاست كه چگونه ميتوان يك ماشين خودكار غير قطعي را بصورت قطعي تبديل نمود. در واقع زيبايي الگوريتمها ها در اين زمينه تبلور خود را نشان مي دهد. براي اين منظور اعداد را در نظر بگيريد . بطور مجزا مي توان گفت كه در حالت كلي سه فرم براي اعداد صحيح و اعشاري مي تواند وجود داشته باشد :



شکل 2.8- انواع اعداد صحيح و اعشاري

در اینجا ترسیم سه ماشین خودکار برای اعداد مستقل از یکدیگر ساده است . مشکل ، ترکیب این سه حالت و ایجاد يك ماشین خودکار برای سه نوع عدد مي باشد. با استفاده از ماشینهاي خودکار غیر قطعي میتوان به سادگی و بدون در نظر گرفتن مشکل ترکیب گراف برای سه نوع متفاوت اعداد اقدام به ترسیم يك ماشین خودکار برای تشخیص سه نوع عدد نمود. زیبایی کار در اینجاست که میتوان بصورت الگوریتمی و بدون نیاز به هیچگونه کار اضافي ، این سه حالت را در قالب يك ماشین خودکار غیر قطعي با یکدیگر ادغام نمود و با استفاده از الگوریتمهایی که ارائه خواهد شد بطور اتوماتیک تبدیل به يك ماشین خودکار قطعي نمود.



شکل 2.9- ماشین خودکار غیر قطعي

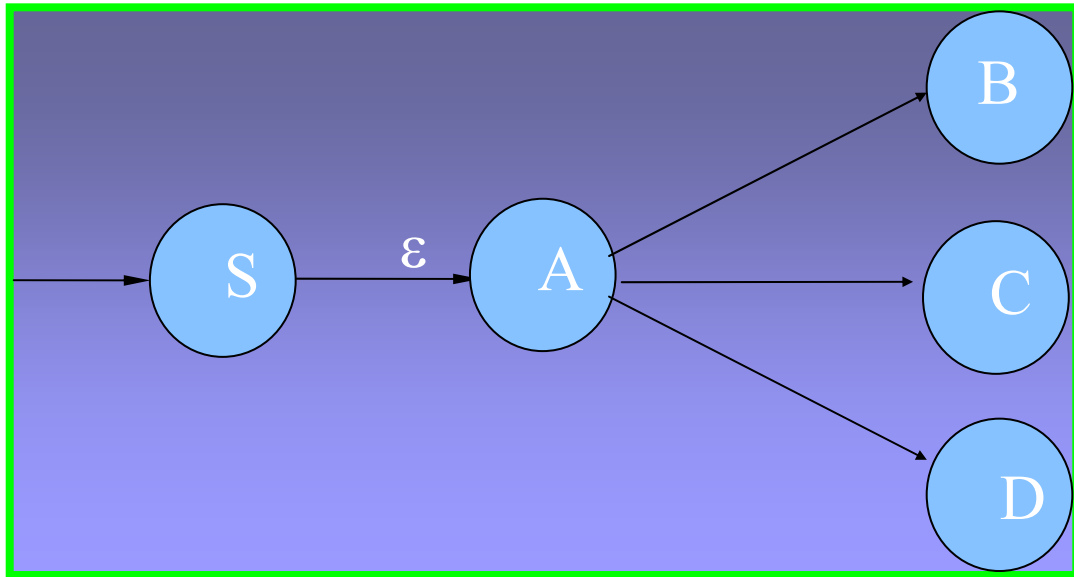
همانگونه که در شکل 2.9 مشاهده مي شود ، يك ماشین خودکار غیر قطعي مي تواند بیان گر حالات مختلف برای يك لغت باشد. در اینجا منظور از لغت ، اعداد مي باشند . اعداد يا داراي علامت هستند و يا علامت ندارند. این نداشتن علامت را با گذر تهی که با علامت ϵ (اسیلون) مشخص شده ، معین مي کنند. همانگونه که مشاهده مي شود، در حالت A با دیدن digit نمی توان مشخص نمود که حالت بعدي آیا حالت B يا C ويا حالت E است . باید توجه داشته باشید که گذر ϵ در واقع یعنی هیچ چیز .

برای تبدیل ماشین خودکار غیر قطعي به قطعي در سه مرحله عمل میشود:

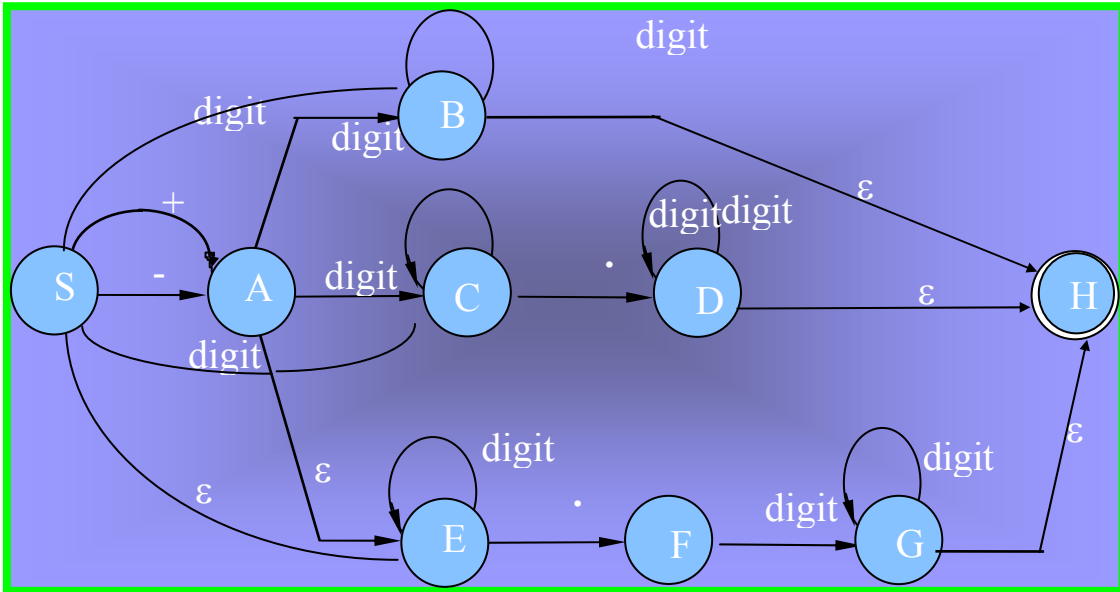
1. حذف گذرهای تهی
2. رفع عدم قطعیت
3. بهینه سازی ماشین خودکار

۲.۹.۱ حذف گذرهای تهی

اگر از حالت S گذری تهی به حالت A وجود دارد، این نمایانگر يك تساوي يك طرفه مي باشد که نکته ايست جالب توجه . در اینجا S با A هیچ تفاوتی ندارد اما A متمایز از S است. میتوان هر واکنشی که در A وجود دارد را برای S نیز در نظر گرفت اما بالعکس صادق نیست.

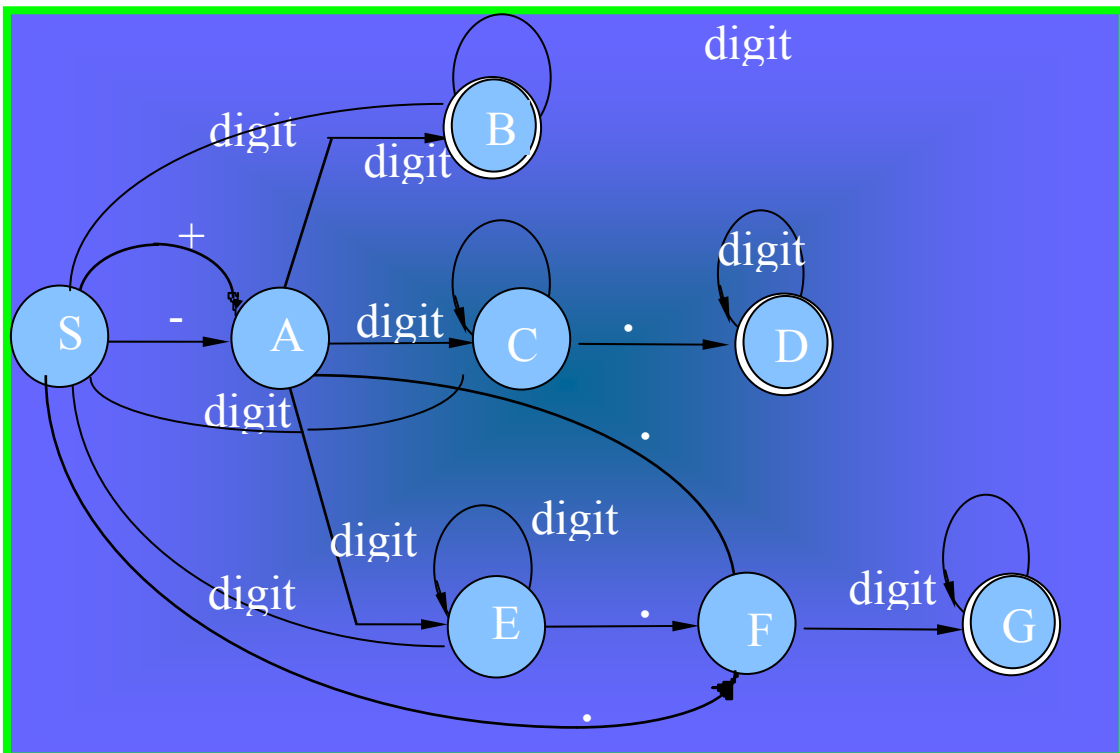


برای درک بهتر مطلب ، گذر تهی از حالت صفر به حالت يك در شکل 2.9 را در نظر بگیرید اگر عدد علامت نداشته باشد، حالت صفر دقیقاً مثل حالت يك است. اما، اگر عدد علامت داشته باشد ، حالت صفر واکنشی دیگر دارد که واکنش این حالت ، گذری به حالت يك می باشد. حالت صفر از حالت يك متمایز گردید ، چرا که ممکن بود دارای علامت باشد، اما اگر عدد علامتی نداشته باشد این دو حالت یکی خواهند بود. پس هیچگونه نیازی به جداسازی آنها نخواهد بود .



شکل 2.10 حذف گذر تهی از حالت صفر به يك

همانگونه که مشاهده می شود، کلیه واکنشهای حالت يك و یا به عبارت دیگر کلیه گذرهای خارج شونده از حالت يك، برای حالت صفر در نظر گرفته شد. به این ترتیب گذر تهی حذف گردید. در مرحله بعد به همین ترتیب ادامه می دهیم تا کلیه گذرهای تهی را از داخل ماشین خودکار قطع نماییم.

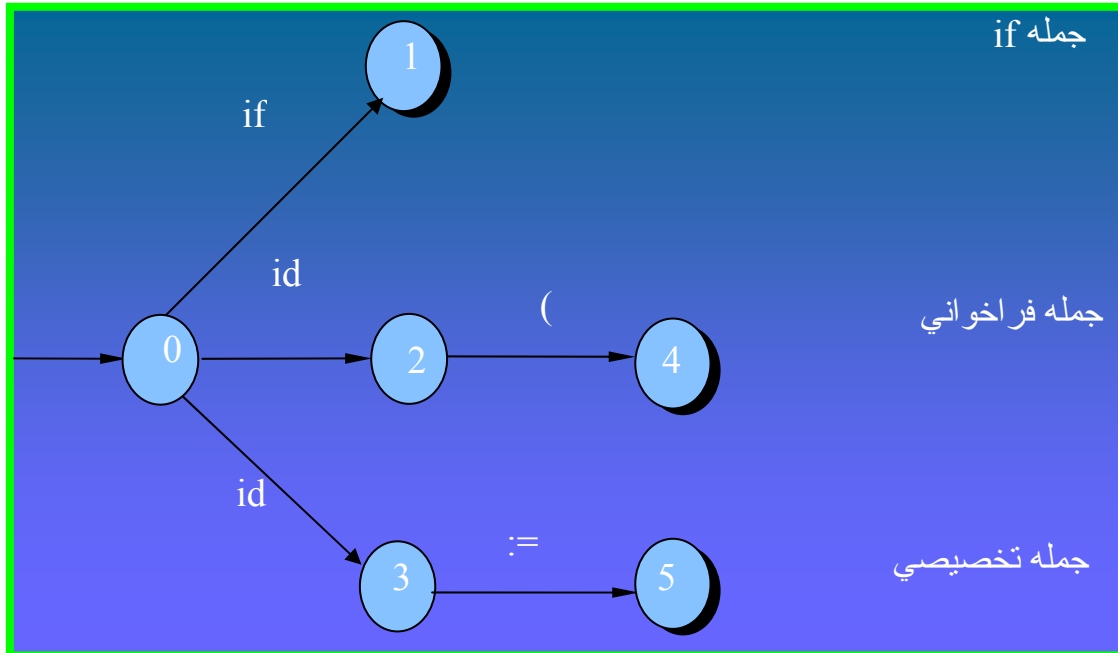


شکل 2.11- ماشین خودکار بدون گذرهای تهی

۲.۹.۲ رفع عدم قطعیت

همانگونه که در شکل 2.11 مشاهده می کنید، ماشین خودکار غیر قطعی است. برای نمونه در حالت A با دیدن digit نمی توان تشخیص داد که آیا حالت بعدی B، C و یا اینکه E می باشد. اما نکته جالب توجه اینجاست که با دیدن digit حتماً به یکی از این سه حالت گذر می شود. برای روشن تر شدن روش حذف عدم قطعیت بهتر است که کار تحلیلگر نحوی مطرح شود.

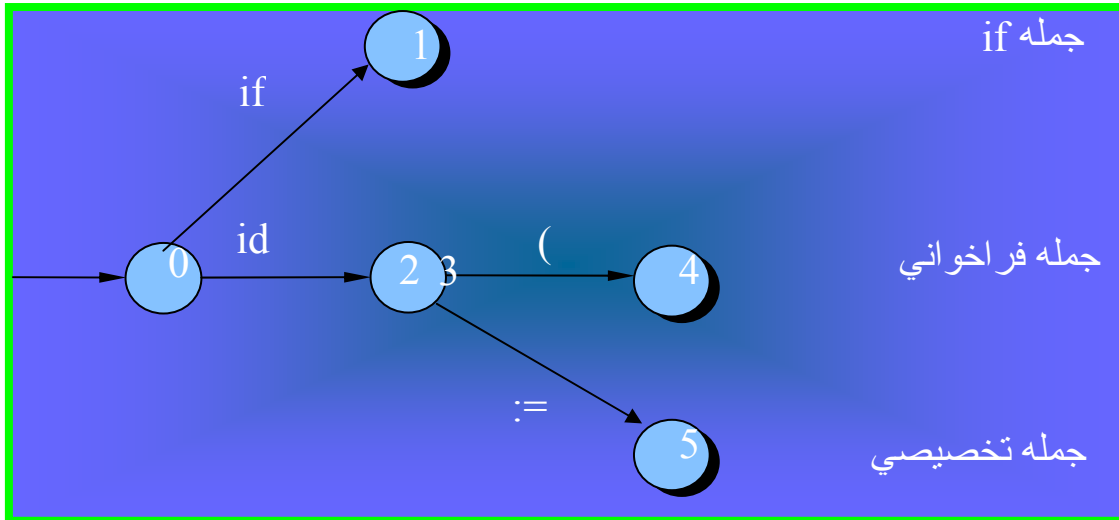
تحلیلگر نحوی، معمولاً با دیدن يك لغت در ورودی، ساختار جمله مورد نظر را می تواند پیش بینی کند. برای مثال اگر لغت دریافتی از تحلیلگر لغوي if باشد، تحلیلگر نحوی بلافاصله تصمیم می گیرد که جمله باید جمله if باشد و باید انتظار دیدن شرط جمله if را در ورودی داشته باشد. اما، با دیدن يك اسم، یا به عبارت دیگر يك شناسه، مشکل وجود دارد. در اینجا تحلیلگر نحوی نمی تواند مشخص کند که آیا جمله مورد نظر يك جمله فراخوانی زیر برنامه ها و یا يك جمله تخصیصی (Assignment) است. ماشین خودکار غیر قطعی در اینجا بصورت زیر قابل ترسیم است:



شکل 2.12 تشخیص جملات if، فراخوانی و تخصیصی

علیرغم وجود عدم قطعیت در حالت صفر از شکل 2.12، می توان قطعیتی را در نظر داشت. به این ترتیب که در حالت شروع صفر اگر يك شناسه یا id در

ورودي ظاهر شود، قطعاً حالت بعدي ، حالت 2 يا 3 خواهد بود. در حالت 2 يا 3 اگر '=' در ورودی ظاهر شود جمله تخصیصی ، اگر '(' ظاهر شود جمله فراخوانی و در غیر اینصورت جمله غلط می باشد. پس حالات 2 يا 3 را بصورت يك حالت ترکیبی جدید با ادغام خروجی های این دو حالت می توان بوجود آورد :



شکل 2.13 ماشین خودکار قطعی

برای سهولت در امر تبدیل ویا نمایش ماشینهای خودکار به جای گراف از فرم جدول مانند استفاده میشود. برای نمونه ماشین خودکار غیر قطعی اعداد در شکل 2.11 ، بصورت يك جدول در شکل 2.14 مشخص شده است.

	digit	.	+/-
S	B,C,E	F	A
A	B,C,E	F	
(B)	B		
C	C	D	
(D)	D		
E	E	F	
F	G		
(G)	G		

شکل 2.14 جدول ماشین خودکار اعداد

همانگونه که در شکل 2.14 مشاهده میشود در حالات S و A به یکی از سه حالت B یا C یا E به ازای دیدن digit گذری وجود دارد. لذا ، میتوان گفت که گذری به حالت ترکیبی BCE وجود دارد. در این حالت میتوان واکنشهای هر سه حالت B ، C و E را داشت. ردیف مربوط به حالت ترکیبی BCE در شکل 2.15 ، ترکیبی از واکنشهای هر سه حالت تشکیل دهنده آن است به عبارت دیگر از ترکیب «یا»

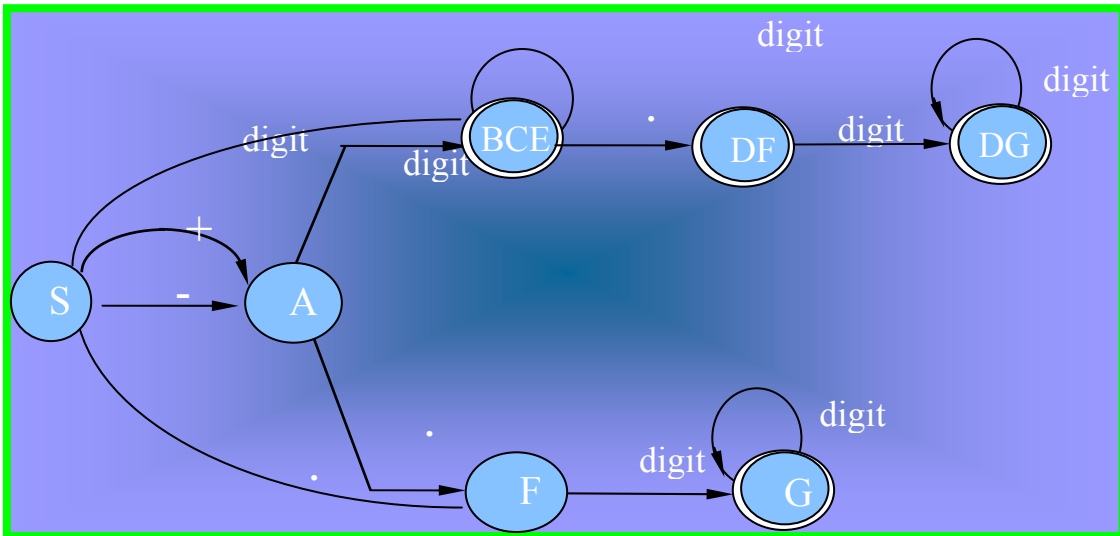
سطرهای B، C و E، سطر جدید BCE به ماتریس افزوده می شود. BCE يك حالت پذیرشي است، زیرا یکی از حالات تشکیل دهنده آن یعنی B يك حالت پذیرش است.

همانگونه که در شکل 2.15 مشاهده میشود، در حالت جدید BCE به ازاء ورودی '.' (نقطه اعشار) عدم قطعیت وجود دارد. از حالت BCE میتوان به هر يك از دو حالت D یا F گذر نمود. لذا برای رفع عدم قطعیت، حالت جدید دیگری بنام DF را به جدول حالات باید افزود. برای رفع عدم قطعیت در حالت DF حالت DG ایجاد میشود.

	digit	.	+/-
S	BCE	F	A
A	BCE	F	
B	B		
C	C	D	
D	D		
E	E	F	
F	G		
G	G		
BCE	BCE	DF	
DF	DG		
DG	DG		

شکل 2.15 جدول ماشین خودکار اعداد

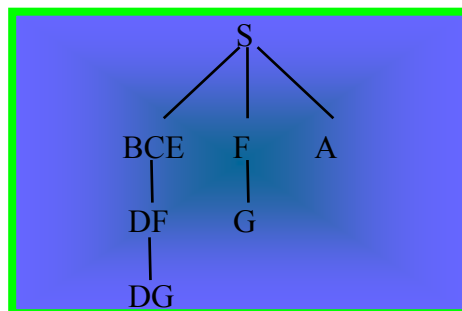
اکنون با استفاده از جدول فوق میتوان ماشین خودکار قطعی اعداد را بصورت زیر ایجاد نمود.



شکل 2.16 ماشین خودکار قطعي اعداد

همانگونه که مشاهده میکنید برخي از حالات درون جدول مثل حالات B ، C ، D و E در عمل غير قابل دسترسي هستند. علت افزايش حالات جديد ترکيبي است. ميتوان اين حالات زائد را با ايجاد درخت دسترسي نیز مشخص نمود.

با افزايش حالات جديد جهت حذف عدم قطعيت در ماشین خودکار ، برخي از حالات زائد شده ، از حالت شروع غير قابل دسترسي ميشوند . جهت تشخيص اين حالات يا مي توان ماشین خودکار را از روي جدول ترسيم نمود و يا اينکه با استفاده از يك درخت دسترسي حالات قابل دسترسي از حالت شروع S را مشخص نمود. براي نمونه درخت دسترسي براي ماتريس شکل 2.15 بصورت زير است. بايد توجه داشته باشيد که نام حالات در درخت دسترسي تکراري نيست و اين درخت صرفاً جهت تعيين حالات قابل دسترسي از حالت شروع S است.



شکل 2.17 درخت حالات قابل دسترسي



هدف از بهینه سازی ، تقلیل تعداد حالات در ماشینهای خودکار است . برای این منظور باید حالات معادل را تشخیص داد . دو حالت را در صورتی معادل گویند که به ازای ورودیهای متفاوت با گذشت از یک سری از حالات یا به عبارت دیگر در مسیری به طول صفر یا بیشتر و با گذشتن از تعداد n گره نهایتاً به حالت پذیرش برسند . حالات معادل را با روشی ابتکاری به ترتیب زیر می توان مشخص نمود :

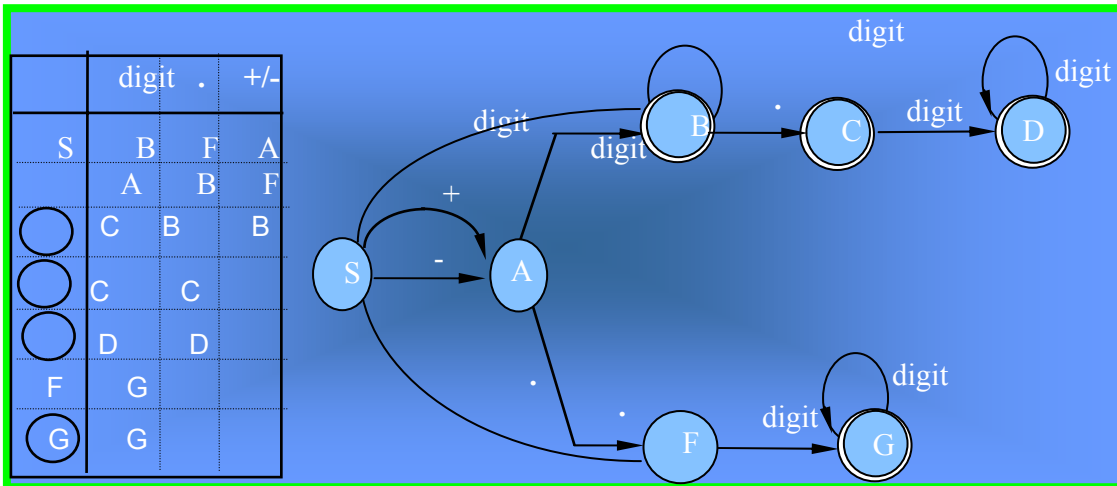
1. ابتدا حالات معادل صفر که با مسیرهایی به طول صفر به پذیرش میرسند را از سایر حالات مجزا باید نمود . بعبارت دیگر ، حالات پذیرش از حالات غیر پذیرش مجزا می شوند . به این ترتیب حالات به دو دسته مجزا از حالات پذیرش و غیر پذیرش تقسیم میشوند .
2. به درون هر دسته حالات باید نگریست و مشخص نمود که آیا ورودی ای وجود دارد که به ازای آن حالات متعلق به یک دسته واکنشهای جدا از سایرین دارند . واکنشها بر اساس گذر یک حالت به حالت دیگر سنجیده نمی شود بلکه ، اگر از یک حالت به حالت دیگر گذری وجود دارد دسته مربوط به آن حالت را مشخص میکنند . مقصود مقایسه دسته هایی است که به آنها گذر میشود . حالات مشابه به ازای یک ورودی به حالات درون یک دسته گذر می کنند .
3. آنقدر مرحله 2 تکرار می شود تا دیگر گذری متفاوت بین دسته ها وجود نداشته باشد و دسته ای جدیدتر تولید نشود . برای نمونه اکنون ماشین خودکار که در مبحث قبل مطرح شد بهینه سازی می شود .

بطور خلاصه و با در نظر گرفتن الگوریتم فوق می توان حالات معادل را بصورت زیر تعریف نمود :

دو حالت را در صورتی معادل k گویند که به ازای هر رشته از ورودیها به طول کوچکتر یا مساوی با k اگر از یکی از آنها بتوان به پذیرش رسید ، از دیگری نیز بتوان به ازاء همان رشته به حالت پذیرش رسید .

دو حالت را معادل گویند اگر و تنها فقط اگر به ازای هر رشته ای از ورودیهای متوالی که موجب رسیدن از آن حالت به یک حالت پذیرش است بتوان از دیگری نیز به پذیرش رسید .

اکنون با توجه به تعریف حالات معادل و مراحل بهینه سازی ماشینهای خودکار ، می توان مبادرت به بهینه سازی ماشین خودکار اعداد که در مبحث قبل مطرح شد ، نمود . ماشین خودکار قطعی اعداد و جدول تولید آن در شکل 2.18 ارائه شده است .



شکل 2.17- ماشین خودکار قطعي اعداد

برای بهینه سازی ماشینهای خودکار ابتدا باید گره ها را به دو دسته پذیرشی و غیر پذیرشی افراز نمود. بعبارت دیگر حالات معادل صفرکه با مسیرهایی به طول صفر به پذیرش می رسند را از سایر حالات باید تفکیک نمود. به این ترتیب حالات ماشین خودکار اعداد که در شکل 2.17 ارائه شده، به دو دسته زیر تقسیم میشوند:

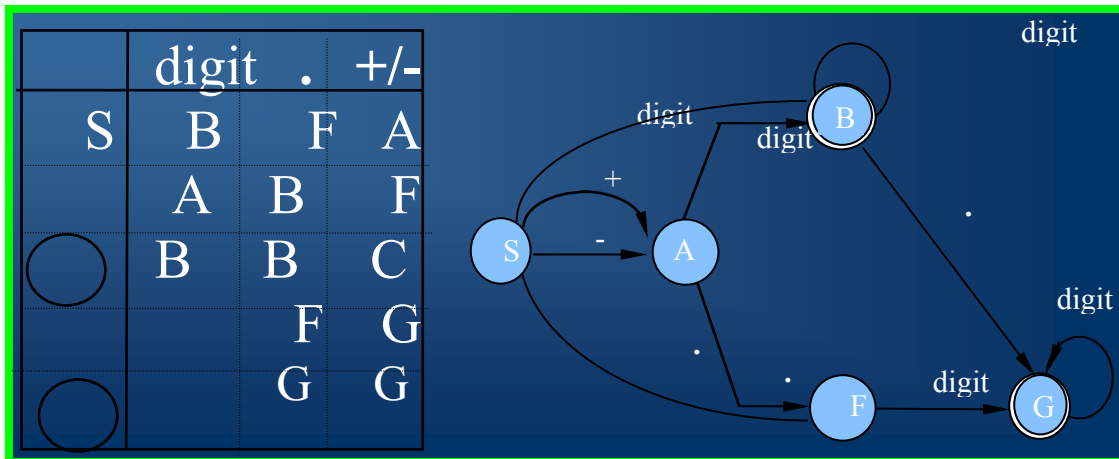
- 1 {B,C,D,G} حالات پذیرشی
- 2 {S,A,F} حالات غیر پذیرشی

حال باید مشخص نمود که آیا در داخل هر دسته به ازاء هر ورودی گره ها واکنشی یکسان دارند. منظور از واکنش یکسان این است که برای نمونه اگر از حالت M در دسته شماره 12 به ازاء ورودی digit گذری به حالتی در دسته شماره 7 وجود دارد سایر حالات موجود در دسته شماره 12 باید به ازاء ورودی digit گذری به حالتی در دسته شماره 7 داشته باشند. در غیر اینصورت این حالات از داخل دسته خارج شده و در دسته ای جدید قرار می گیرند. برای نمونه در داخل دسته شماره 1 در بالا به ازاء ورودی نقطه اعشار '،' از حالت B يك خروجی و گذری وجود دارد در صورتیکه در مورد سایر حالات اینچنین نیست. پس این حالت از سایرین جدا میشود و سه دسته به صورت زیر حاصل می گردد:

- 1 { C ,D ,G}
- 2 (B)
- 3 {S ,A ,F} حالات غیر پذیرشی

به همین ترتیب در دسته شماره 3 حالت S قابل تفکیک از A و F است. زیرا در S به ازاء ورودیهای + و - گذر و واکنشی وجود دارد در صورتیکه در مورد A و F

اینچنین نیست. پس حالت S از A و F قابل تفکیک است. از سوي دیگر A و F نیز قابل تفکیک از یکدیگر هستند. بنابراین سه حالت C، D و G غیر قابل تفکیک از یکدیگر و معادل هستند. پس میتوان هر يك از این سه حالت معادل را بجای دو حالت دیگر در داخل ماشین خودکار قرار داد. به این ترتیب ماشین خودکار بهینه اعداد بصورت زیر خواهد بود.

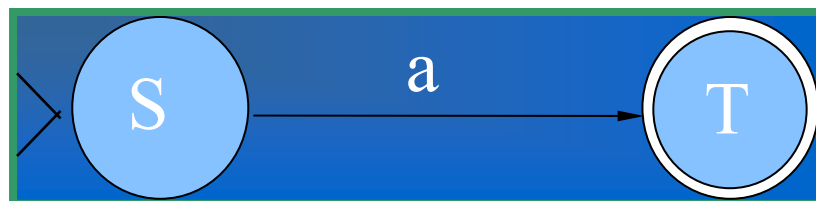


شکل 2.18- ماشین خودکار بهینه اعداد

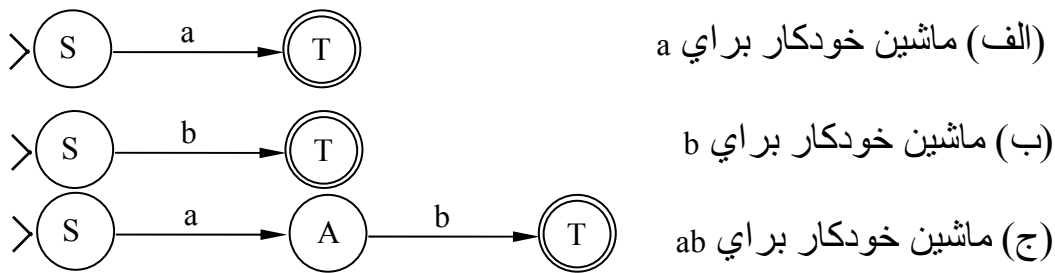
۲.۱۰ تبدیل عبارات با قاعده به ماشینهای خودکار

در این قسمت نشان داده خواهد شد که چگونه می توان عبارات با قاعده را به صورت ماشینهای خودکار تبدیل نمود تا اینکه بتوان با استفاده از ماشین خودکار، قوانین لغوي که در فرم عبارات باقاعده خلاصه شده اند را عملاً به صورت کد برنامه مورد بهره برداري قرار داد. برای تبدیل عبارات با قاعده به ماشینهای خودکار معمولاً مراحل زیر بکار می روند:

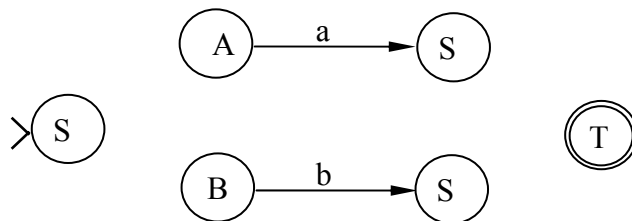
1. هر عنصر a متعلق به الفبای زبان به صورت يك ماشین خودکار نمایش داده می شود:



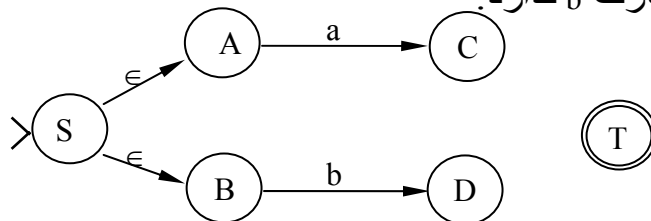
2. ماشین خودکار برای عبارت ab را از ترکیب گرافها برای a و b به صورت زیر می توان ایجاد نمود.



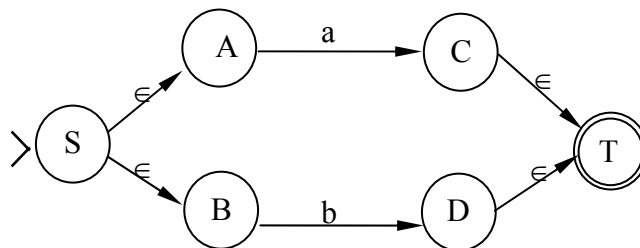
3. ماشین خودکار برای عبارت $a|b$ را با يك استدلال ساده میتوان از ماشین های خودکار برای عبارات a و b ایجاد نمود. چنانچه حالت شروع ماشین خودکار برای عبارت $a|b$ حالت S باشد:



حالت شروع S هیچ فرقی با حالت شروع برای عبارت a ندارد زیرا در شروع $a|b$ میتوان a را هم در ورودی دید. از سوی دیگر حالت شروع S هیچ تفاوتی با حالت شروع برای عبارت b ندارد.



به همین ترتیب با مشاهده a در ورودی باید مطمئن بود که $a|b$ در ورودی ظاهر شده است. لذا، حالت پذیرش برای ماشین خودکار عبارت a یعنی C هیئت تفاوتی با حالت پذیرش برای ماشین خودکار برای عبارت $a|b$ یعنی T ندارد. به همین ترتیب حالت D هیچ تفاوتی با حالت پذیرش T ندارد. بنابراین ماشین خودکار برای عبارت $a|b$ بصورت زیر خواهد بود:

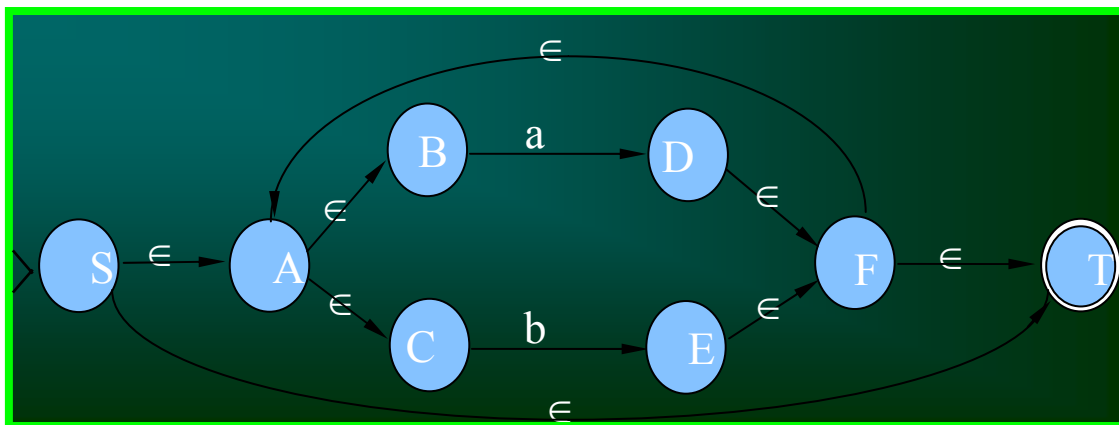




در شکل فوق عبارت $a|b$ از ترکیب ماشینهای خودکار برای عبارات a و b ساخته شده است. مسلم است که حالت شروع S هیچ تفاوتی با حالت A و حالت B ندارد. اما بالعکس صادق نیست و حالت شروع برای تشخیص a یعنی A متفاوت از حالت شروع برای b است زیرا در حالت شروع برای a نمی توان در ورودی b را دید. لذا، حالت شروع $a|b$ را با دو گذر تهی به حالت شروع برای a و حالت شروع برای b متصل باید نمود.

4. با در دست داشتن ماشین خودکار برای a و b و در نتیجه وجود ماشین خودکار برای $a|b$ می توان ماشین خودکار برای عبارت $(a|b)^*$ را با یک استدلال ساده ایجاد نمود. به این ترتیب که فرض کنید حالات شروع و خاتمه برای ماشین خودکار $(a|b)^*$ به ترتیب حالات S و T باشند. اولاً، چون $(a|b)^*$ میتواند تهی نیز باشد پس حالت شروع آن ممکن است هیچ تفاوتی با حالت پذیرش نداشته باشد. زیرا رشته $(a|b)^*$ ممکن است اصلاً وجود نداشته باشد.

از طرف دیگر ممکن است حالت شروع S هیئ تفاوتی با حالت شروع $(a|b)$ نداشته باشد زیرا، در شروع $(a|b)^*$ میتوان در شروع مشاهده $a|b$ در ورودی بود. پس از مشاهده $a|b$ دو باره میتوان در آغاز مشاهده $a|b$ جدیدتری قرار گرفت. این در واقع بخاطر خاصیت تکراری بودن $(a|b)^*$ است. پس، حالت پذیرش $a|b$ درون $(a|b)^*$ ، هیچ تفاوتی با حالت شروع آن نخواهد داشت. با این استدلال میتوان نتیجه گرفت که ماشین خودکار برای عبارت $(a|b)^*$ بصورت زیر است.

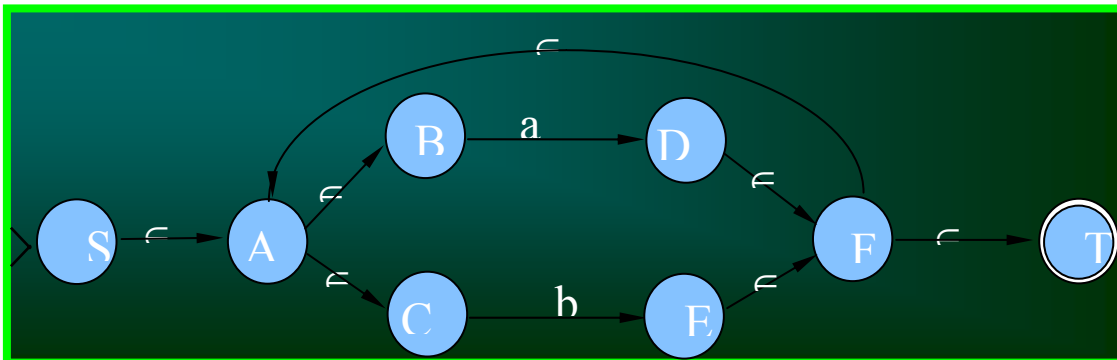


شکل 2.19- ماشین خودکار برای عبارت $(a|b)^*$

مسلماً در ورودی $(a|b)^*$ نون رشته a یا b می تواند اصلاً وجود نداشته باشد با این نتیجه حالت شروع می تواند معادل با حالت پذیرش باشد. لذا، در شکل 2.19 حالت شروع S با گذری تهی به حالت پذیرش T متصل شده است. از جهت دیگر در خاتمه دیدن $b|a$ مثل اینکه دوباره در حالت شروع بوده و می توان $b|a$

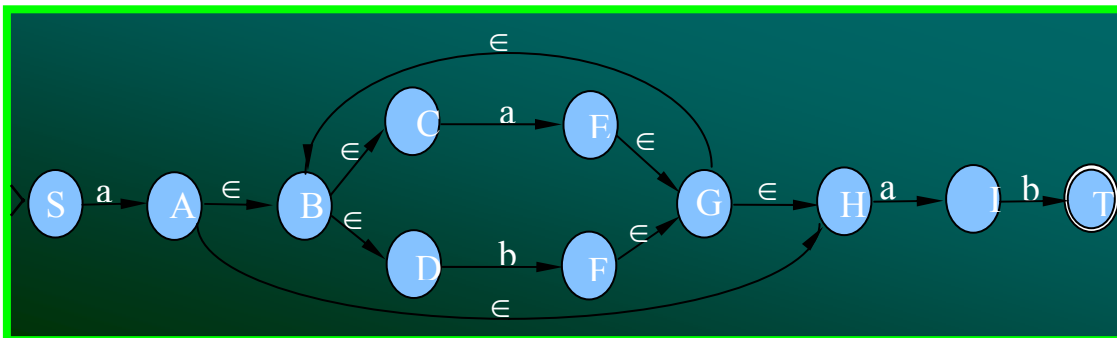
جديدي را مشاهده نمود و اين عمل تا به هر تعدادي قابل تکرار است . براي اين منظور حالت F با گذري تهی به حالت شروع $a|b$ يعني A متصل شده است. علاوه بر اين حالت F ممکن است خاتمه تکرار و حالت پذیرش هم باشد.

5 . عبارت $(a|b)^+$ را مي توان با حذف گذر تهی از حالت شروع به حالت پايان در ماشين خودکار براي $(a|b)^*$ توليد نمود زیرا در مورد $(a|b)^+$ حداقل یکبار $a|b$ در ورودی باید ظاهر شود و حالت شروع آن با خاتمه یکسان نیست .



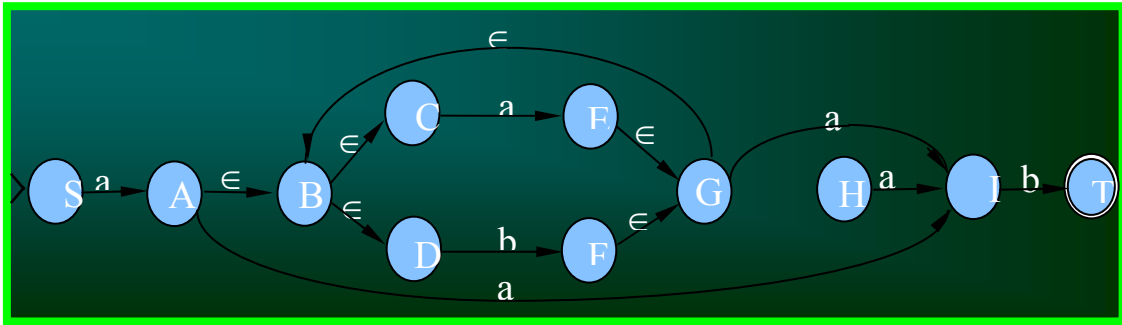
شکل 2.20- ماشين خودکار براي عبارت $(a|b)^+$

مثال 1- براي عبارت $a(a|b)^*a$ يك ماشين خودکار بهينه ايجاد نماييد. با استفاده از ماشين خودکار ارائه شده براي عبارات $(a|b)^*$ که در شکل 2.19 ارائه شده است ، ميتوان بسادگی ماشين خودکار غير قطعي را ايجاد کرد.



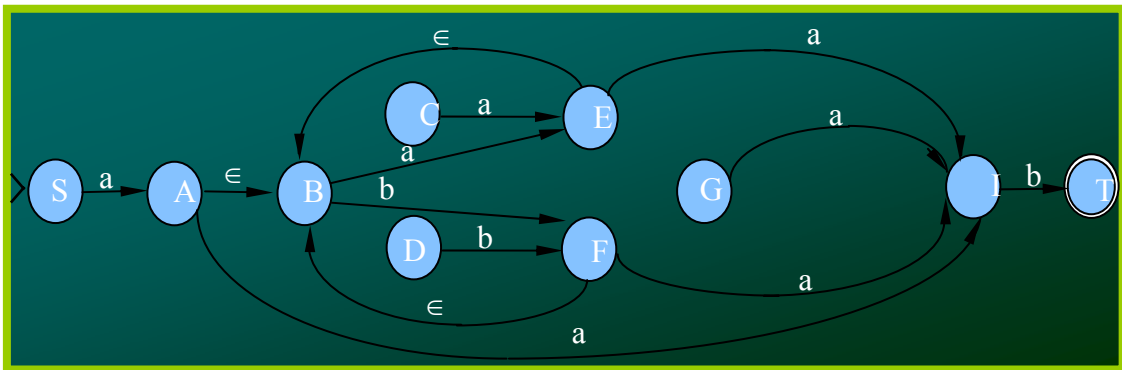
شکل 2.21- ماشين خودکار غير قطعي براي عبارت

پس از حذف گذر هاي تهی به H ماشين بصورت زیر تبديل ميشود.



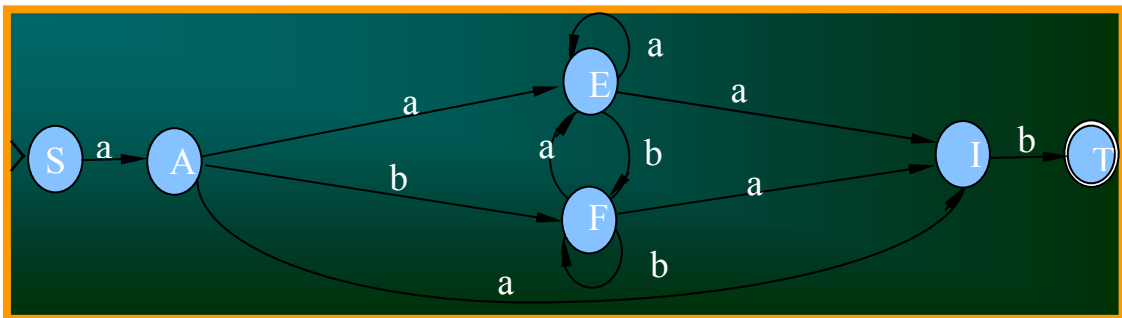
شکل 2.22- حذف گذرهای تهی از A و G به H

مشاهده میکنید که حالت H در شکل 2.22 غیر قابل دسترسی و قابل حذف است. در شکل 2.23 گذرهای تهی به حالات G، C و D حذف شده اند.



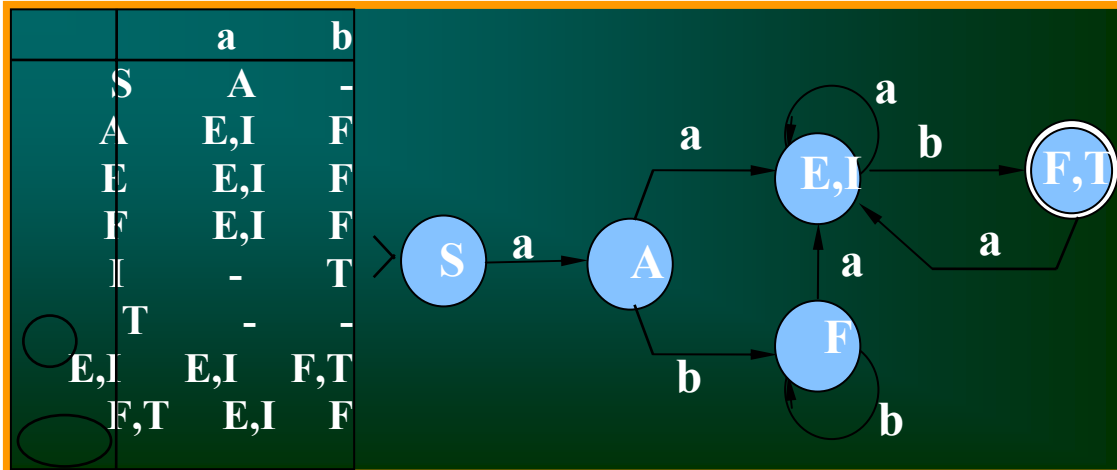
شکل 2.23- حذف گذرهای تهی به حالات G، C و D

پس از حذف گذرهای تهی ماشین بصورت زیر تبدیل میشود.



شکل 2.24- ماشین خودکار پس از حذف گذرهای تهی

ماشین خودکار فوق در شکل 2.25 با استفاده از نمایش ماتریسی به فرم قطعی تبدیل شده است.



شکل 2.25- ماشین خودکار قطعی برای عبارت $a(a|b)^*ab$

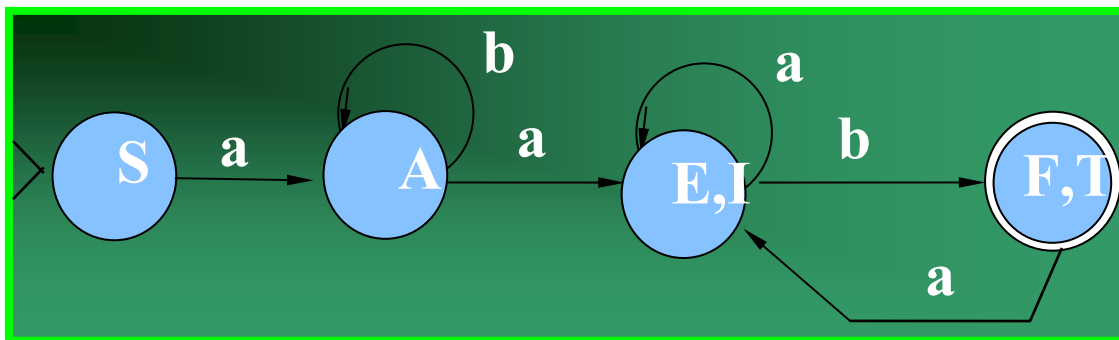
برای بهینه سازی ماشین خودکار فوق حالات مربوطه را به دو دسته پذیرشی و غیر پذیرشی بصورت زیر میتوان تقسیم بندی نمود:

الف) 1- {FT} 2- {S, A, EI, F}

ب) 1- {FT} 2- {S} 3- {A, EI, F}

ج) 1- {FT} 2- {S} 3- {A, F} 4- {EI}

به این ترتیب مشاهده میکنید که حالات A و F با یکدیگر معادل هستند و یکی را میتوان با دیگری جایگزین نمود. به این ترتیب ماشین خودکار بهینه برای عبارت $a(a|b)^*ab$ بصورت زیر خواهد بود.





۲.۱۰ تمرین

تمرین 1- فرم کلي اعدادي را مشخص کنید که یا صفرمي باشندویا اینکه اگر طولشان بیش از يك است حتماً "با يك رقم غیر صفر آغاز می شوند.

تمرین 2- فرم کلي اعداد در مبنای هشت را مشخص کنید در صورتی که بدانیم اینگونه اعداد حتماً با يك رقم صفر آغاز و سپس حرف 0 ظاهر می گردد و بعد از آن عدد مبنای هشت مشخص می شود. برای نمونه 00127 يك عدد مبنای 8 است.

تمرین 3- فرم کلي اعداد در مبنای شانزده در زبان C با رقم صفر و سپس X آغاز می شود. يك عبارت با قاعده برای بیان اعداد در مبنای شانزده ارائه نمایید. بطور مثال 0XAF25 يك عدد در مبنای 16 است.

تمرین 4- فرم کلي رشته ها را در زبان C را با يك عبارت با قاعده معین کنید. توجه داشته باشید که در زبان C علامت \ (Backslash) در داخل يك رشته مفهوم خاص دارد.

تمرین 5- فرم کلي اعداد صحیحی را تعیین کنید که در آنها ارقام به ترتیب صعودی ظاهر میشوند.

تمرین 6- فرم کلي کلیه رشته اعداد مبنای دو را مشخص کنید که در آنها زیر رشته 011 ظاهر نگردد.

تمرین 7 * - کلیه رشته های اعداد مبنای دو را مشخص کنید که تعداد صفرها در آنها فرد و تعداد يك ها زوج باشد.

تمرین 8- ماشین خودکار قطعی برای شناسه ها ترسیم نمائید

تمرین 9- ماشین خودکار قطعی بهینه برای عبارت $(a|b|c)^*abb$ ایجاد نماید و سپس برنامه تحلیلگر لغوي برای تشخیص اینگونه رشته ها بنویسید. ابتدا عبارت را به فرم غیر قطعی تبدیل نمائید سپس قطعی و بهینه نمائید.

تمرین 10- يك ماشین خودکار قطعی برای کلیه اعداد مبنای دو که تعداد صفرهای آنها فرد و تعداد يك ها فرد است ایجاد کنید. در ضمن لااقل دو عدد صفر و دو عدد يك باید در این اعداد موجود باشد.



تمرین 11- ماشین خودکار قطعي براي اعداد مبناي 2 که حتما" داراي يك رقم صفر و يك رقم يك هستند را ايجاد نمائيد بقسمي که اعداد با تعداد يك ها و صفرهاي زوج را در يك حالت پذيرشي مشخص نمايد . اعداد با تعداد يك هاي فرد و تعداد صفرهاي فرد را نیز در حالت پذيرش ديگر بايد ماشین مشخص کند. براي اعداد با تعداد يك ها فرد و صفرها زوج و همچنين صفرها فرد و يك ها زوج نیز باید دو حالت پذيرش مجزا وجود داشته باشد. ابتدا براي هر يك از چهار حالت ذکر شده ماشینهاي خودکار را ترسيم نمائيد و سپس يك ماشین خودکار غير قطعي براي ترکیب آنها ايجاد نمائيد. سپس این ماشین خودکار را قطعي کنید.

تمرین 12- يك ماشین خودکار قطعي براي گرامر يکي از زبانهاي C يا پاسکال ايجاد نمائيد. سپس با استفاده از روش ارائه شده در این فصل يك تحلیلگر لغوي براي آن ايجاد کنید.