

/home/ted/awk_stuff/awk.language.summary

Prepared, on the basis of *The AWK Programming Language* by Aho, Kernighan & Weinberger,
by Ted Harding (original 01/12/2001; mods: 20/08/2002, 07/10/2004, 24/08/2006, 07/06/2007, 17/07/2007,
19/04/2008, 09/09/2010)

The AWK Crib

Basic statement

pattern { *action* }

or *function definition*

Pattern Summary

Summary of Patterns

1. **BEGIN** { *statements* }
The *statements* are executed once, before any input has been read.
BEGIN does not combine with any other pattern, and **does** require an action.
 2. **END** { *statements* }
The *statements* are executed once, after all input has been read.
END does not combine with any other pattern, and **does** require an action.
 3. *expression* { *statements* }
The *statements* are executed at each input line where the *expression* is true (non-zero or non-null).
 4. */regular expression/* { *statements* }
The *statements* are executed at each input line matched by the *regular expression*.
!/regular expression/ { *statements* }
The *statements* are executed at each input line **not** matched by the *regular expression*.
 5. *compound pattern* { *statements* }
A *compound pattern* combines expressions with **&&** (AND), **||** (OR), **!** (NOT) and parentheses; the *statements* are executed at each input line where the *compound pattern* is true.
 6. *pattern₁ , pattern₂* { *statements* }
A **range pattern** matches each input line, starting from a line matched by *pattern₁* and ending with the next line matched by *pattern₂* (both inclusive); the *statements* are executed at each such line.
A **range pattern** can not be part of any other pattern.
 7. { *statements* }
The **missing pattern** matches every input line (including blank lines), so the *statements* will be executed for every line.
 8. *pattern*
A pattern with no explicit { *action* } has an implicit action which simply outputs the whole of any input line matched by *pattern*.
-
-

Expressions as Patterns

Any *expression* can be used as a *pattern*. If it evaluates to **non-zero** or **non-null** for a line, then the line matches and any actions are executed.

Expressions and Operators

Any *expression* can be used as **operand** for any **operator**.

Expressions and Operators

Expressions

- The primary expressions are
 - numeric and string constants;
 - variables;
 - fields;
 - function calls;
 - array elements.

Operators

- The following operators combine expressions

the assignment operators	= += -= *= /= %= ^=
the conditional expressions operator	? : <i>cond</i> ? <i>expr</i> ₁ : <i>expr</i> ₂
the logical operators	&& !
the substring matching operators	~ !~
the relational operators	< <= == != >= >
the string concatenation operator	blank
the arithmetic operators	+ - * / % ^
unary arithmetic	+ -
prefix and suffix increment and decrement operators	++ --
parentheses for grouping	

String-Matching Patterns

String-Matching Patterns

- /regexpr/*
Matches whenever the current line contains a **substring** matched by *regexpr*.
 - !/regexpr/*
Matches whenever the current line **does not** contain a **substring** matched by *regexpr*.
 - expression ~ regexpr*
expression ~ expression₂
Matches whenever the string value of *expression* contains a substring matched by *regexpr*, or equal to the string value of *expression₂*.
 - expression !~ regexpr*
expression !~ expression₂
Matches whenever the string value of *expression* does not contain a substring matched by *regexpr*, or equal to the string value of *expression₂*.
-

Regular Expressions

Metacharacters

\ ^ . [] | () * + ?

Expression	Matches
<i>c</i>	the non-metacharacter <i>c</i>
<i>\c</i>	escape sequence, or literal <i>c</i> \b backspace \f formfeed \n newline \r carriage return \t tab \ddd octal value <i>ddd</i> \c any other character <i>c</i> literally
<i>^</i>	beginning of string
<i>\$</i>	end of string
<i>.</i>	any character
<i>[c₁c₂...]</i>	any character in <i>c₁c₂...</i>
<i>[^c₁c₂...]</i>	any character not in <i>c₁c₂...</i>
<i>[c₁ - c₂]</i>	any character in the range <i>c₁ - c₂</i>
<i>[^c₁ - c₂]</i>	any character not in the range <i>c₁ - c₂</i>
<i>r₁ r₂</i>	any string matched by <i>r₁</i> or by <i>r₂</i>
<i>(r₁)(r₂)</i>	any string <i>xy</i> where <i>r₁</i> matches <i>x</i> and <i>r₂</i> matches <i>y</i>
<i>(r)*</i>	zero or more consecutive strings matched by <i>r</i>
<i>(r)+</i>	one or more consecutive strings matched by <i>r</i>
<i>(r)?</i>	zero or one string matched by <i>r</i>
<i>(r)</i>	any string matched by <i>r</i>

Actions

An *action* consists of one or more *statements* separated by newlines or semicolons.

Action Statements

1. *expression*
2. **print** *expression-list*
print(*expression-list*)
3. **printf** *format, expression-list*
printf(*format, expression-list*)
4. **if** (*expression*) *statement*
5. **if** (*expression*) *statement* **else** *statement* [**else** associated with most recent unmatched **if**]
6. **while** (*expression*) *statement*
7. **for** (*expression; expression; expression*) *statement*
8. **for** (*variable in array*) *statement*
9. **do** *statement* **while** (*expression*)
10. **break**
11. **continue**
12. **next**
13. **exit**
14. **exit** *expression*
15. { *statements* }

User-defined Variables

Named user-defined variables come into existence on being mentioned for the first time. They have default null initial values of "" (in a string context) or 0 (in a numerical context). Variables can also be defined in the command-line.

Variables acquire values by assignment. The value of a string variable is automatically converted to numeric type (if this makes sense[†]) if the variable is used in a numerical context. The value of a numeric variable is automatically converted to string type[‡] if the variable is used in a string context.

[†] If not, the numeric value used is 0, the default numeric null initial value.

E.g. `{i="hello";printf("%5.5f ",i); print i}` → 0.00000 hello

The numeric value of a string is the value of its longest initial numerically compatible sequence,

e.g. "0.12abc" → "0.12" → 0.12, "abc" → "" → 0, "1E2G3H4" → "1E2" → 100

[‡] The string value is its representation according to OFMT (default "% .6g"),

e.g. 100/3 → "33.3333", 100000000/3 → "3.33333e+07".

Functions use **local** copies of variables named in their argument lists. All other variables (including variables defined within functions) are **global**; variables arising within a function which are intended to be local can be so coerced by being named (with default null values) in the argument list of the function definition.

Built-in Variables

Variable	Meaning	Default
ARGC	number of command-line arguments	
ARGV	array of command-line arguments	
ARGV[i]	ith command-line argument	
FILENAME	name of current input file	
FNR	current record number in current file	
FS	field separator (character or string)	" "‡
NF	number of fields in current record	
NR	number of records read so far	
OFMT	output format for numbers	"%.6g" (i.e. 6 significant figures)
OFS	output field separator	" "
ORS	output record separator	"\n"
RLENGTH	length of string matched by <code>match()</code>	
RS	input record separator	"\n"
RSTART	start of string matched by <code>match()</code>	
SUBSEP	subscript separator (<code>ctrl-\ = ^\</code>) [†]	"\034"

[†] For the **awk** 'hack' of multidimensional arrays
(See the A-K-W book pp. 52–3, and below)

[‡] Multiple `FS=" ..."` treated as one " ";
for `FS="[1[1..."`, each " " is significant.
String `FS` is treated as regular expression
(leftmost longest non-overlapping match).
`FS="[1[1*|[\t]|[:]|[,]"` will cause
one or more spaces, or **TAB**, or ":", or ",",
to be recognised as a field separator.
`FS` can be a string, e.g. "(" is set up by
`FS="()|(|)"`

Field Variables

Variable	Meaning	Default
\$0	The whole line	
\$1, \$2, ...	The fields of the line	
\$(i)	Field i (dynamic index)	
\$var	E.g. <code>{Fno=25 ; print(\$Fno)}</code> (dynamic index)	
\$(NF+1) = ...	Defines an additional field; increments NF	

Functions

There are several **built-in** functions. The user may define any number of **user-defined** functions.

User-defined Functions

Such a function is defined by a *function definition* statement (declaration) of the form

```
function name ( parameter-list ) { statements }
```

which may occur anywhere a *pattern-action* statement can. See above for the distinction between **local** and **global** variables. Example:

```
function log10(x) { return log(x)/log(10) }
```

In calling a function there must be no space between the name and the opening parenthesis of the parameter-list (a space is permissible in the definition).

Functions may be defined recursively. Internal variables used in recursive functions need to appear in the parameter list of the function declaration (so that they will be *local* at every depth; otherwise they are *global* and the function will misbehave).

Built-in Arithmetic Functions

FUNCTION	VALUE RETURNED
atan2 (<i>y</i> , <i>x</i>)	arctangent of <i>y</i> / <i>x</i> in the range $-\pi$ to π (radians)
cos (<i>x</i>)	cosine of <i>x</i> (<i>x</i> in radians)
exp (<i>x</i>)	exponential function of <i>x</i> , e^x
int (<i>x</i>)	integer part of <i>x</i> ; truncated towards 0 when $x > 0$
log (<i>x</i>)	natural (base <i>e</i>) logarithm of <i>x</i>
rand ()	random number [†] <i>r</i> , uniform real on $0 \leq r < 1$
sin (<i>x</i>)	sine of <i>x</i> (<i>x</i> in radians)
sqrt (<i>x</i>)	square root of <i>x</i>
srand (<i>x</i>)	<i>x</i> (unsigned integer) is new seed [†] for rand ()

[†] **rand**() starts, by default, from the same value (1) every time **awk** is run.

srand(*x*) starts **rand**() from the value *x*; **srand**() sets *x* from the system clock.

srand returns the *current* (**pre-call**) seed; to get (for re-use) the seed set by **srand**() do

```
srand(); i=srand(); print i; srand(i);
```

Built-in String Functions

FUNCTION	ACTION, AND VALUE RETURNED
gsub (<i>r</i> , <i>s</i>)	substitute <i>s</i> for regexp <i>r</i> globally in \$0 ; <i>return</i> number of substitutions made
gsub (<i>r</i> , <i>s</i> , <i>t</i>)	substitute <i>s</i> for regexp <i>r</i> globally in string <i>t</i> ; <i>return</i> number of substitutions made
index (<i>s</i> , <i>t</i>)	<i>return</i> first position of string <i>t</i> in <i>s</i> , or 0 if not present
length (<i>s</i>)	<i>return</i> number of characters in <i>s</i>
match (<i>s</i> , <i>r</i>)	test whether <i>s</i> contains a substring matched by regexp <i>r</i> ; <i>return</i> index or 0; sets RSTART and RLENGTH (leftmost longest [†]) <i>e.g.</i> match(\$0,/[0-9]+\ .htm/)
split (<i>s</i> , <i>A</i>)	split <i>s</i> into array <i>A</i> on FS ; [‡] <i>return</i> number of fields
split (<i>s</i> , <i>A</i> , <i>fs</i>)	split <i>s</i> into array <i>A</i> on field separator <i>fs</i> ; [‡] <i>return</i> number of fields
sprintf (<i>fmt</i> , <i>expr-list</i>)	<i>return</i> <i>expr-list</i> formatted according to <i>fmt</i>
sub (<i>r</i> , <i>s</i>)	substitute <i>s</i> for leftmost longest [†] substring of \$0 matched by regexp <i>r</i> ; <i>return</i> number of substitutions made
sub (<i>r</i> , <i>s</i> , <i>t</i>)	substitute <i>s</i> for leftmost longest [†] substring of <i>t</i> matched by regexp <i>r</i> ; <i>return</i> number of substitutions made
substr (<i>s</i> , <i>p</i>)	<i>return</i> last part of <i>s</i> starting at position <i>p</i>
substr (<i>s</i> , <i>p</i> , <i>n</i>)	<i>return</i> substring of <i>s</i> of length <i>n</i> starting at position <i>p</i>

[†] First the leftmost match is found; then it is extended as far as possible

[‡] Index values in **A** are integers 1, 2, ...; Element values **A**[*i*] are substrings of *s*

Operation

It is important to note that **awk** *always* operates by applying its program to each input line, and generating appropriate output each time. If no input lines are supplied, nothing will happen. For example, let `{...}` denote the above **awk** program for fish names. Then the command line

```
awk '{...}'
```

will apparently hang. In fact, **awk** is waiting for input lines and will generate appropriate output for any lines which are matched (as, in this case, all are). Therefore, nothing will happen unless **Return** is pressed, but the above output will be generated every time it is pressed; press **^D** (**EOF**) to stop it.