

**Confidentiality Statement**

This document contains confidential information of Tata Consultancy Services Limited, which is provided for the sole purpose of permitting the recipient to evaluate the proposal submitted herewith. In consideration of receipt of this document, the recipient agrees to maintain such information in confidence and to not reproduce or otherwise disclose this information to any person outside the group directly responsible for evaluation of its contents, except that there is no obligation to maintain the confidentiality of any information which was known to the recipient prior to receipt of such information from Tata Consultancy Services Limited, or becomes publicly known through no fault of recipient, or is received without obligation of confidentiality from a third party owing no obligation of confidentiality to Tata Consultancy Services Limited.

### **Tata Code of Conduct**

We, in our dealings, are self-regulated by a Code of Conduct as enshrined in the Tata Code of Conduct. We request your support in helping us adhere to the Code in letter and spirit. We request that any violation or potential violation of the Code by any person be promptly brought to the notice of the Local Ethics Counselor or the Principal Ethics Counselor or the CEO of TCS. All communication received in this regard will be treated and kept as confidential.



## **TATA CONSULTANCY SERVICES**

### **Virtual ILP – Introduction to awk Filter**

---

#### **Content Manual**

Version 1.1

**December 2014**  
(ILP Guwahati)

## Table of Content

|  |    |
|--|----|
| 1.AWK inbuilt variables.....           | 6  |
| 1.1.FS (Input Field Separator).....    | 7  |
| 1.2.OFS (Output Field Separator) ..... | 8  |
| 1.3.RS (Row separator).....            | 8  |
| 1.4.ORS (Output Record Separator)..... | 8  |
| 1.5.NF.....                            | 9  |
| 1.6.NR .....                           | 10 |
| 1.7.FNR.....                           | 10 |
| 1.8.FILENAME.....                      | 11 |
| 2.Awk Built in Function.....           | 11 |
| 2.1.Arithmetic Functions.....          | 11 |
| 2.2. String Functions.....             | 14 |

## 1. *AWK inbuilt variables*

AWK is supplied with good number of built-in variables which comes in handy when working with data files. We will see each AWK built-in variables with one or two examples to familiarize with them. Without these built-in variables it's very much difficult to write simple AWK code. These variable are used to format output of an AWK command, as input field separator and even we can store current input file name in them for using them within the script.

|                 |  |
|-----------------|--|
| <b>FS</b>       | field separator character (default blank & tab)    |
| <b>OFS</b>      | output field separator string (default blank)      |
| <b>RS</b>       | input record separator character (default newline) |
| <b>ORS</b>      | output record separator string (default newline)   |
| <b>NF</b>       | number of fields in input record                   |
| <b>NR</b>       | number of input record                             |
| <b>FNR</b>      | output number of lines                             |
| <b>FILENAME</b> | name of current input file                         |

Consider below db.txt as sample file.

**~\$ cat db.txt**

```
John,29,MS,IBM,M,Married
Barbi,45,MD,JHH,F,Single
Mitch,33,BS,BofA,M,Single
Tim,39,Phd,DELL,M,Married
Lisa,22,BS,SmartDrive,F,Married
```

In order to make it simple we can divide above inbuilt variables in to groups on basis of their operations.

**Group1:** FS(input field separator), OFS

**Group2:** RS(Row separator) and ORS(Output record separator)

**Group3:** NR, NF and FNR

**Group4:** FILENAME variable

### 1.1. *FS (Input Field Separator)*

This variable is useful in storing the input field separator. By default AWK can understand only spaces, tabs as input and output separators. But if your file contains some other character as separator other than these mentioned ones, AWK cannot understand them.

For example [UNIX password file](#) which contains ':' as a separator. So in order to mention the input field separator we use this inbuilt variable. We will see what issue we face if we don't mention the field separator for our db.txt.

**Example: without using FS**

Print first column data from db.txt file.

```
~$ awk '{print $1}' db.txt
```

**Output:**

```
John,29,MS,IBM,M,Married
Barbi,45,MD,JHH,F,Single
Mitch,33,BS,BofA,M,Single
Tim,39,Phd,DELL,M,Married
Lisa,22,BS,SmartDrive,F,Married
```

If you see entire file is displayed which indicates AWK does not understand db.txt file separator ','. We have to tell AWK what is the field separator.

**Example: Using FS**

List only first column data from db.txt file which has field separator as ','.

```
~$ awk 'BEGIN{FS=","}{print $1}' db.txt
```

**Output:**

```
John
Barbi
Mitch
Tim
Lisa
```

## 1.2. *OFS (Output Field Separator)*

This variable is useful for mentioning what is your output field separator which separates output data.

**Example:**

Display only 1st and 4th column and the separator between at output for these columns should be \$.

```
~$ awk 'BEGIN{FS=",";OFS=" $ "}{print $1,$4}' db.txt
```

**Output:**

John \$ IBM

Barbi \$ JHH

Mitch \$ BofA

Tim \$ DELL

Lisa \$ SmartDrive

**Note:** Space is give before and after \$ in OFS variable to show better output.

## 1.3. *RS (Row separator)*

Row Separator is helpful in defining separator between rows in a file. By default AWK takes row separator as new line. We can change this by using RS built-in variable.

**Example:**

Convert a sentence to a word per line. We can use RS variable for doing it.

```
~$ echo "This is how it works" | awk 'BEGIN{RS=" "} {print $0}'
```

**Output:**

This

is

how

it

works

### **1.4.      *ORS (Output Record Separator)***

This variable is useful for defining the record separator for the AWK command output. By default ORS is set to new line.

**Example:**

Print all the company names in single line which are in 4th column.

```
~$ awk -F',' 'BEGIN{ORS=" "} {print $4}' db.txt
```

**Output:**

IBM JHH BofA DELL SmartDrive

### **1.5.      *NF***

This variable keeps information about total fields in a given row. The final value of a row can be represented with \$NF.

**Example:** Consider abc.txt which contains below data:

Jones 2143 78 84 77  
Gondrol 2321 56 58 45  
RinRao 2122234 38 37  
Edwin 253734 87 97 95  
Dayan 24155 30 47

Print number of fields in each row in abc.txt.

```
~$ awk '{print NF}' abc.txt
```

**Output:**

5  
5  
4  
5  
4



**Example:**

Print last field in each row of abc.txt file.

***~\$ awk '{print \$NF}' abc.txt***

**Output:**

77

45

37

95

47

**Note:** If you observe above two examples we used Just NF for giving us the count of fields in a given row and \$NF for displaying last element in each row. \$NF will come handy when you are not sure what your last column number is.

## **1.6. NR**

This variable keeps the value of present line number. This will come handy when you want to print line numbers in a file.

**Example:**

Print line number for each line in a given file.

***~\$ awk '{print NR, \$0}' abc.txt***

**Output:**

1 Jones 2143 78 84 77

2 Gondrol 2321 56 58 45

3 RinRao 2122234 38 37

4 Edwin 253734 87 97 95

5 Dayan 24155 30 47

This can be treated as [cat command -n](#) option for displaying line number for a file.

### **1.7. *FNR***

This variable keeps count of number of lines present in a given file/data. This will come handy when you want to print no of line present in a given file. This command is equivalent to wc -l command.

**Example:**

Print total number of lines in a given file.

```
~$ awk 'END{print FNR}' abc.txt
```

**Output:**

5

### **1.8. *FILENAME***

This variable contain file awk command is processing.

**Example:**

Print filename for each line in a given file.

```
~$ awk '{print FILENAME, NR, $0}' abc.txt
```

**Output:**

```
abc.txt 1 Jones 2143 78 84 77
abc.txt 2 Gondrol 2321 56 58 45
abc.txt 3 RinRao 2122234 38 37
abc.txt 4 Edwin 253734 87 97 95
abc.txt 5 Dayan 24155 30 47
```

## **2. *Awk Built in Function***

A function is a self-contained computation that accepts a number of arguments as input and returns some value. Awk has a number of built-in functions in two groups: arithmetic and string functions.

## 2.1. Arithmetic Functions

Nine of the built-in functions can be classified as arithmetic functions. Most of them take a numeric argument and return a numeric value. Below table summarizes these arithmetic functions with some examples.

| Awk Function                               | Description   |
|--|---|
| <code>cos ( <i>x</i> )</code>              | Returns cosine of <i>x</i> ( <i>x</i> is in radians).   |
| <code>exp ( <i>x</i> )</code>              | Returns e to the power <i>x</i> .   |
| <code>index (s1,s2)</code>                 | Position of string s2 in s1; returns 0 if not present   |
| <code>int ( <i>x</i> )</code>              | Returns truncated value of <i>x</i> .   |
| <code>log ( <i>x</i> )</code>              | Returns natural logarithm (base- <i>e</i> ) of <i>x</i> .   |
| <code>sin ( <i>x</i> )</code>              | Returns sine of <i>x</i> ( <i>x</i> is in radians).   |
| <code>sqrt ( <i>x</i> )</code>             | Returns square root of <i>x</i> .   |
| <code>atan2 ( <i>y</i> , <i>x</i> )</code> | Returns arctangent of <i>y</i> / <i>x</i> in the range $-\pi$ to $\pi$ .                          |
| <code>rand ()</code>                       | Returns pseudo-random number <i>r</i> , where $0 \leq r < 1$ .                                    |
| <code>srand ( <i>x</i> )</code>            | Establishes new seed for rand(). If no seed is specified, uses time of day. Returns the old seed. |

### Examples:

`sqrt(expr)`

```
~$ awk 'BEGIN{
print sqrt(16);
print sqrt(0);
print sqrt(-12);
}'
```

### Output:

```
4
0
nan
Here nan stands for not a valid number.
```

`sin(expr)`

```
~$ awk 'BEGIN {  
  print sin(90);  
  print sin(45);  
}'
```

**Output:**

```
0.893997  
0.850904
```

**index()**

```
~$ awk 'BEGIN { print Index("peanut", "an") }'
```

**Output:**

```
3
```

**substr(s,m,n)**

```
~$ awk '{ print substr("deepak",2,3) }'
```

**Output**

```
ee
```

**rand()**

```
~$ cat rand.awk  
BEGIN {  
  while(i<1000)  
  {  
    n = int(rand()*100);  
    rnd[n]++;  
    i++;  
  }  
  for(i=0;i<=100;i++) {  
    print i,"Occurred", rnd[i], "times";  
  }  
}
```

```
}
}
```

```
~$awk -f rand.awk
```

**Output:**

```
0 Occurred 6 times
1 Occurred 16 times
2 Occurred 12 times
3 Occurred 6 times
4 Occurred 13 times
.
.
.
```

## 2.2. *String Functions*

The built-in string functions are much more significant and interesting than the numeric functions. Because awk is essentially designed as a string-processing language, a lot of its power derives from these functions. [Below](#) table lists the string functions found in awk.

### Awk's Built-In String Functions

| Awk Function  | Description  |
|---|--|
| <code>gsub ( <i>r</i> , <i>s</i> , <i>t</i> )</code>    | Globally substitute's <i>s</i> for each match of the regular expression <i>r</i> in the string <i>t</i> . Returns the number of substitutions. If <i>t</i> is not supplied, defaults to <b>\$0</b> . |
| <code>index ( <i>s</i> , <i>t</i> )</code>              | Returns position of substring <i>t</i> in string <i>s</i> or zero if not present.  |
| <code>length ( <i>s</i> )</code>                        | Returns length of string <i>s</i> or length of <b>\$0</b> if no string is supplied.  |
| <code>match ( <i>s</i> , <i>r</i> )</code>              | Returns either the position in <i>s</i> where the regular expression <i>r</i> begins or 0 if no occurrences are found. Sets the values of <b>RSTART</b> and <b>RLENGTH</b> .                         |
| <code>split ( <i>s</i> , <i>a</i> , <i>sep</i> )</code> | Parses string <i>s</i> into elements of array <i>a</i> using field separator <i>sep</i> ; returns  |

| Awk Function                          | Description   |
|---------------------------------------|---|
|                                       | number of elements. If <i>sep</i> is not supplied, <b>FS</b> is used. Array splitting works the same way as field splitting.<br>Uses <code>printf</code> format specification for <b>expr</b> . |
| <code>sprintf ( "fmt ", expr )</code> |   |
| <code>sub ( r , s , t )</code>        | Substitutes <i>s</i> for first match of the regular expression <i>r</i> in the string <i>t</i> . Returns 1 if successful; 0 otherwise. If <i>t</i> is not supplied, defaults to <b>\$0</b> .    |
| <code>substr ( s , p , n )</code>     | Returns substring of string <i>s</i> at beginning position <i>p</i> up to a maximum length of <i>n</i> . If <i>n</i> is not supplied, the rest of the string from <i>p</i> is used.             |
| <code>tolower ( s )</code>            | Translates all uppercase characters in string <i>s</i> to lowercase and returns the new string.   |
| <code>toupper ( s )</code>            | Translates all lowercase characters in string <i>s</i> to uppercase and returns the new string.   |