

DATA PROCESSING WITH AWK

A COMMAND-LINE TOOL TO SLICE & DICE TEXTUAL DATA

KRISHNAKUMAR GOPALAKRISHNAN

krishnak [at] vt [dot] edu

COMPUTATIONAL METHODS HUB

IMPERIAL COLLEGE LONDON

APRIL 16 2019

Imperial College
London

INTRODUCTION

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al
- Weinberger, Peter

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al
- Weinberger, Peter
- Kernighan, Brian

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al
- Weinberger, Peter
- Kernighan, Brian

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al
- Weinberger, Peter
- Kernighan, Brian

What can we do with Awk?

- Text manipulation in powerful ways

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al
- Weinberger, Peter
- Kernighan, Brian

What can we do with Awk?

- Text manipulation in powerful ways
- Data filtering, cleaning & other pre-processing tasks

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al
- Weinberger, Peter
- Kernighan, Brian

What can we do with Awk?

- Text manipulation in powerful ways
- Data filtering, cleaning & other pre-processing tasks
- Read & write data in a variety of formats

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al
- Weinberger, Peter
- Kernighan, Brian

What can we do with Awk?

- Text manipulation in powerful ways
- Data filtering, cleaning & other pre-processing tasks
- Read & write data in a variety of formats
- Produce tabular reports for the web (advanced)

WHAT IS AWK?

Why the *awkward* name?

- Aho, Al
- Weinberger, Peter
- Kernighan, Brian

What can we do with Awk?

- Text manipulation in powerful ways
- Data filtering, cleaning & other pre-processing tasks
- Read & write data in a variety of formats
- Produce tabular reports for the web (advanced)
- Task automation (advanced)

WHERE IS AWK AVAILABLE?

Summary of Awk Variants

- Lots of variants (awk, bawk, nawk, mawk, gawk, goawk)
- Most popular version of awk is GNU Awk (gawk)
- Version 5.0 released Apr 12, 2019!

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants
 - ▶ Solaris, Illumos & many others

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants
 - ▶ Solaris, Illumos & many others

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants
 - ▶ Solaris, Illumos & many others

Windows Machines

- Plethora of options available
 - ▶ Cygwin: complete unix-like environment (heavy)

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants
 - ▶ Solaris, Illumos & many others

Windows Machines

- Plethora of options available
 - ▶ Cygwin: complete unix-like environment (heavy)
 - ▶ Windows Subsystem for Linux (WSL): Windows 10 1703 & above (requires admin privileges)

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants
 - ▶ Solaris, Illumos & many others

Windows Machines

- Plethora of options available
 - ▶ Cygwin: complete unix-like environment (heavy)
 - ▶ Windows Subsystem for Linux (WSL): Windows 10 1703 & above (requires admin privileges)
 - ▶ ezwinports (<https://sourceforge.net/projects/ezwinports/>)

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants
 - ▶ Solaris, Illumos & many others

Windows Machines

- Plethora of options available
 - ▶ Cygwin: complete unix-like environment (heavy)
 - ▶ Windows Subsystem for Linux (WSL): Windows 10 1703 & above (requires admin privileges)
 - ▶ ezwinports (<https://sourceforge.net/projects/ezwinports/>)
 - ▶ Git for Windows (<https://git-scm.com/download/win>)

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants
 - ▶ Solaris, Illumos & many others

Windows Machines

- Plethora of options available
 - ▶ Cygwin: complete unix-like environment (heavy)
 - ▶ Windows Subsystem for Linux (WSL): Windows 10 1703 & above (requires admin privileges)
 - ▶ ezwinports (<https://sourceforge.net/projects/ezwinports/>)
 - ▶ Git for Windows (<https://git-scm.com/download/win>)
 - ▶ Cmder (<https://github.com/cmderdev/cmder/releases/download/v1.3.11/cmder.zip>)

WHERE IS AWK AVAILABLE?

*nix Machines

- Pre-installed on most Unix-like OSes
 - ▶ MacOS
 - ▶ Linux
 - ▶ All BSD variants
 - ▶ Solaris, Illumos & many others

Windows Machines

- Plethora of options available
 - ▶ Cygwin: complete unix-like environment (heavy)
 - ▶ Windows Subsystem for Linux (WSL): Windows 10 1703 & above (requires admin privileges)
 - ▶ ezwinports (<https://sourceforge.net/projects/ezwinports/>)
 - ▶ Git for Windows (<https://git-scm.com/download/win>)
 - ▶ Cmder (<https://github.com/cmderdev/cmder/releases/download/v1.3.11/cmder.zip>)
 - ▶ Log on to a *nix remote server

APPLICABILITY OF AWK

- Great for:

APPLICABILITY OF AWK

- Great for:
 - ▶ Manipulating text files divided into lines and columns

APPLICABILITY OF AWK

- Great for:
 - ▶ Manipulating text files divided into lines and columns
 - ▶ All lines are not required to be in the same format

APPLICABILITY OF AWK

- Great for:

- ▶ Manipulating text files divided into lines and columns
- ▶ All lines are not required to be in the same format
- ▶ Performs best on a structured file (eg tabular data)

APPLICABILITY OF AWK

- Great for:
 - ▶ Manipulating text files divided into lines and columns
 - ▶ All lines are not required to be in the same format
 - ▶ Performs best on a structured file (eg tabular data)
- Small one-line Awk programs can:

APPLICABILITY OF AWK

- Great for:
 - ▶ Manipulating text files divided into lines and columns
 - ▶ All lines are not required to be in the same format
 - ▶ Performs best on a structured file (eg tabular data)
- Small one-line Awk programs can:
 - ▶ Find interesting lines in a data file

APPLICABILITY OF AWK

- Great for:
 - ▶ Manipulating text files divided into lines and columns
 - ▶ All lines are not required to be in the same format
 - ▶ Performs best on a structured file (eg tabular data)
- Small one-line Awk programs can:
 - ▶ Find interesting lines in a data file
 - ▶ Output only columns of data matching some criterion

APPLICABILITY OF AWK

- Great for:
 - ▶ Manipulating text files divided into lines and columns
 - ▶ All lines are not required to be in the same format
 - ▶ Performs best on a structured file (eg tabular data)
- Small one-line Awk programs can:
 - ▶ Find interesting lines in a data file
 - ▶ Output only columns of data matching some criterion
 - ▶ Swapping the order of columns

APPLICABILITY OF AWK

- Great for:
 - ▶ Manipulating text files divided into lines and columns
 - ▶ All lines are not required to be in the same format
 - ▶ Performs best on a structured file (eg tabular data)
- Small one-line Awk programs can:
 - ▶ Find interesting lines in a data file
 - ▶ Output only columns of data matching some criterion
 - ▶ Swapping the order of columns
 - ▶ Combine multiple columns into one

APPLICABILITY OF AWK

■ Great for:

- ▶ Manipulating text files divided into lines and columns
- ▶ All lines are not required to be in the same format
- ▶ Performs best on a structured file (eg tabular data)

■ Small one-line Awk programs can:

- ▶ Find interesting lines in a data file
- ▶ Output only columns of data matching some criterion
- ▶ Swapping the order of columns
- ▶ Combine multiple columns into one
- ▶ Split single column into multiple

APPLICABILITY OF AWK

■ Great for:

- ▶ Manipulating text files divided into lines and columns
- ▶ All lines are not required to be in the same format
- ▶ Performs best on a structured file (eg tabular data)

■ Small one-line Awk programs can:

- ▶ Find interesting lines in a data file
- ▶ Output only columns of data matching some criterion
- ▶ Swapping the order of columns
- ▶ Combine multiple columns into one
- ▶ Split single column into multiple
- ▶ Sophisticated data operations (joins, merges etc)

APPLICABILITY OF AWK

■ Great for:

- ▶ Manipulating text files divided into lines and columns
- ▶ All lines are not required to be in the same format
- ▶ Performs best on a structured file (eg tabular data)

■ Small one-line Awk programs can:

- ▶ Find interesting lines in a data file
- ▶ Output only columns of data matching some criterion
- ▶ Swapping the order of columns
- ▶ Combine multiple columns into one
- ▶ Split single column into multiple
- ▶ Sophisticated data operations (joins, merges etc)

■ Not good for:

APPLICABILITY OF AWK

■ Great for:

- ▶ Manipulating text files divided into lines and columns
- ▶ All lines are not required to be in the same format
- ▶ Performs best on a structured file (eg tabular data)

■ Small one-line Awk programs can:

- ▶ Find interesting lines in a data file
- ▶ Output only columns of data matching some criterion
- ▶ Swapping the order of columns
- ▶ Combine multiple columns into one
- ▶ Split single column into multiple
- ▶ Sophisticated data operations (joins, merges etc)

■ Not good for:

- ▶ Manipulating binary files (Excel/Word)

APPLICABILITY OF AWK

■ Great for:

- ▶ Manipulating text files divided into lines and columns
- ▶ All lines are not required to be in the same format
- ▶ Performs best on a structured file (eg tabular data)

■ Small one-line Awk programs can:

- ▶ Find interesting lines in a data file
- ▶ Output only columns of data matching some criterion
- ▶ Swapping the order of columns
- ▶ Combine multiple columns into one
- ▶ Split single column into multiple
- ▶ Sophisticated data operations (joins, merges etc)

■ Not good for:

- ▶ Manipulating binary files (Excel/Word)
- ▶ Not a web programming language (eg parsing HTML)

HANDS-ON EXERCISES

EXERCISE 1

- Reverse names in `names.txt`

QUESTIONS?