

## Handy Command-line One-liners for Starting Data Scientists

(from <https://medium.com/@rama100/handy-command-line-one-liners-for-starting-data-scientists-81f933241128>)

Category	Task	One-Liner	Comments
File exploration	Count the number of lines in a file	<code>wc -l filename</code>	
	Show the column names, one in each line, preceded by line numbers (i.e., grab the header row, transpose it and prefix line numbers)*	<code>head -1 filename   tr '\t' '\n'   nl</code> <code>head -1 filename   tr ';' '\n'   nl</code> <code>head -1 filename   tr ' ' '\n'   nl</code>	Helpful when you have numerous columns in a new file and want to get the lay of the land e.g., knowing that "average_selling_price" is column # 39 is useful for some of the examples below.
	Page through the file, showing line numbers	<code>less filename   nl</code>	
Row-oriented operations	Show line #4212	<code>sed '4212q;d' filename</code>	Very useful when you are trying to load the file into a database and the load fails at line #4212, for instance. Also, this command will conveniently quit after printing the 4212nd line; very considerate if your file has a million lines!
	Show lines with "foo" in any field	<code>grep 'foo' filename</code>	
	show lines with "foo" in any field, ignoring case	<code>grep -i 'foo' filename</code>	
	Show lines with 'foo' in field #18*	<code>awk -F\t '\$18 == "foo"' filename</code> <code>awk -F, '\$18 == "foo"' filename</code> <code>awk '\$18 == "foo"' filename</code>	
	Show rows that have fewer fields than the header row	<code>awk 'NR==1 {x=NF}; NF &lt; x' filename</code>	Quick check to see if any rows are incomplete
	Remove lines with "foo" in any field and save the rest into a new file	<code>sed '/foo/d' filename &gt; newfile</code>	
	Remove lines with 'foo' in field #18 and save the rest into a new file*	<code>awk -F\t '\$18 != "foo"' filename &gt; newfile</code> <code>awk -F, '\$18 != "foo"' filename &gt; newfile</code> <code>awk '\$18 != "foo"' filename &gt; newfile</code>	
	Remove the first line and save the rest into a new file	<code>sed '1d' filename &gt; newfile</code>	Great for stripping a header row before further processing
	Remove the first 8 lines and save the rest into a new file	<code>sed '1,8d' filename &gt; newfile</code>	
	Remove line #42 and save the rest into a new file	<code>sed '42d' filename &gt; newfile</code>	
Remove lines 233 to 718 and save the rest into a new file	<code>sed '233,718d' filename &gt; newfile</code>		

\* The first one-liner is for tab-delimited files, the second for comma-delimited files, and the third for space-delimited files.

	Remove the last line and save the rest into a new file	<code>sed '\$d' filename &gt; newfile</code>	
	Remove the last 8 lines and save the rest into a new file	<code>sed -e :a -e '\$d;N;2,8ba' -e 'P;D' filename &gt; newfile</code>	Ugh! Let me know if you know of a better way
	Remove blank lines from the file and save the rest into a new file	<code>sed '/^\$/d' filename &gt; newfile</code>	
	Remove duplicate lines and save the rest into a new file, preserving the original order	<code>awk '!seen[\$0]++' filename &gt; newfile</code>	
	Remove duplicate lines and save the rest into a new file, original order may not be preserved	<code>sort -u filename &gt; newfile</code>	
	Remove lines with a missing value in field #18 and save the rest into a new file*	<code>awk -F\t '!\$18' filename &gt; newfile</code> <code>awk -F, '!\$18' filename &gt; newfile</code> <code>awk '!\$18' filename &gt; newfile</code>	
Column-oriented operations	Show just col #42*	<code>cut -f42 filename</code> <code>cut -d, -f42 filename</code> <code>cut -d' ' -f42 filename</code>	
	Show the unique values in column #42 with counts*	<code>cut -f42 filename   sort   uniq -c</code> <code>cut -d, -f42 filename   sort   uniq -c</code> <code>cut -d' ' -f42 filename   sort   uniq -c</code>	Useful for understanding a categorical field. A histogram, essentially.
	Remove the 1st field and save the rest into a new file*	<code>cut -f2- filename &gt; newfile</code> <code>cut -d, -f1-18,43- filename &gt; newfile</code> <code>cut -d' ' -f1-18,43- filename &gt; newfile</code>	
	Remove field #42 and save the rest into a new file*	<code>cut -f1-41,43- filename &gt; newfile</code> <code>cut -d, -f1-41,43- filename &gt; newfile</code> <code>cut -d' ' -f1-41,43- filename &gt; newfile</code>	
	Remove fields #19-42 and save the rest into a new file*	<code>cut -f1-18,43- filename &gt; newfile</code> <code>cut -d, -f1-18,43- filename &gt; newfile</code> <code>cut -d' ' -f1-18,43- filename &gt; newfile</code>	
Combining/splitting files	Stack files column-wise	<code>paste file1 file2 &gt; newfile</code>	Useful if you have two or more files with the same rows but different sets of columns and you need to combine them side-by-side
	Stack files row-wise	<code>cat file1 file2 &gt; newfile</code>	Useful if you have two or more files with the same columns and you need to 'pancake' stack them. Assumes file2 doesn't have a header row. If it does, first remove it using a one-liner :-)
	Split a file into two files with 3000 rows in the first file and the rest in the second	<code>csplit -sf prefix filename 3001</code>	Useful for train/test splitting (the resulting files will be prefixed with whatever you specify as 'prefix')

\* The first one-liner is for tab-delimited files, the second for comma-delimited files, and the third for space-delimited files.

Random selection	Randomly shuffle the rows and save to a new file	<code>awk 'BEGIN{srand();}{print rand()"\t"\$0}' filename   sort -k1 -n   cut -f2- &gt; newfile</code>	
	Randomly choose X% of rows and save to a new file (X= 10% in the code snippet)	<code>awk -v X=10 'BEGIN {srand()} rand() &lt;= 0.01*X' filename &gt; newfile</code>	
	Select every 10th row and save to a new file	<code>awk 'NR%10' filename &gt; newfile</code>	

\* The first one-liner is for tab-delimited files, the second for comma-delimited files, and the third for space-delimited files.