CS 124 / LINGUIST 180

Stanford University

Text Processing with UNIX

Adapted from Ken Church's "Unix for Poets" by Dan Jurafsky and Chris Manning



Text processing with UNIX

Billions of words are text are everywhere

The web, email, social media

What can we do with it all?

The UNIX intuition:

- Modularity and simplicity! Simple command-line tools combined together
- Often faster even than writing a quick python tool
- DIY is very satisfying

7 exercises we'll be doing today

- 1. Count words in a text with "tr", "sort", and "uniq"
- 2. Counting more complex text patterns with "tr"
- 3. Using "sort" in more powerful ways for counting
- 4. Using "grep" to find patterns
- 5. Data ethics: Where did this data come from?
- 6. Simple tricks to compute n-gram statistics
- 7. Using "sed" to modify text



Tools

sort **uniq** –**c** (count duplicates) **tr** (translate characters) **grep**: search for a pattern (regular expression) **cat** (send file(s) in stream) head tail **rev** (reverse lines) **sed** (edit string -- replacement)

Prereqs:

Mac: Open the Terminal app

Windows:

Make sure you have Ubuntu (PAO)

Prerequisites: get the text file we are using

- rice: ssh into a farmshare and then do (don't forget the final) ".")
- cp /afs/ir/class/cs124/WWW/nyt 200811.txt
- Or download to your own laptop this file: http://cs124.stanford.edu/nyt 200811.txt Or:
- scp rice:/afs/ir/class/cs124/WWW/nyt 200811.txt

•

Prerequisites

The UNIX "man" command

- e.g., man tr
- man shows you the command options
 - it's not particularly friendly

Prerequisites

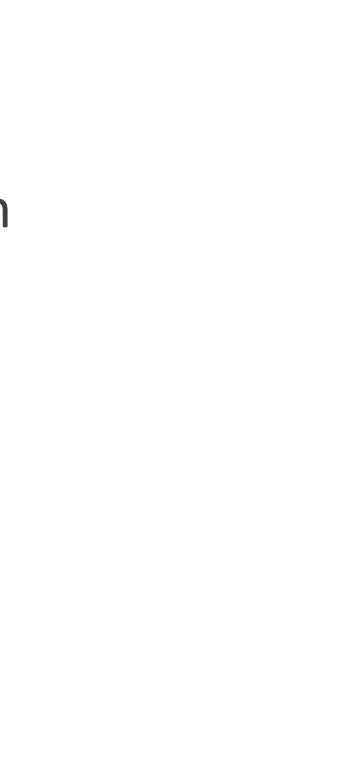
How to chain shell commands and deal with input/output

Input/output redirection:

- > "output to a file"
- < "input from a file"
- | "pipe"

CTRL-C

The less command (quit by typing "q")



Exercise 1: Count words in a text

Input: text file (nyt 200811.txt)

Output: list of words in the file with freq counts Algorithm

- 1. Tokenize (tr)
- 2. Sort (sort)
- **3. Count duplicates (**uniq -c)

Go read the man pages and figure out how to pipe these together

Solution to Exercise 1

• tr -sc 'A-Za-z' '\n' < nyt 200811.txt | sort | uniq -c

633 A

- 1 AA
- AARP

ABBY 1

41 ABC

1 ABCNews

(Do you get a different sort order?) In some versions of UNIX, sort doesn't use ASCII order (uppercase before lowercase).)

Some of the output

- tr -sc 'A-Za-z' '\n' < nyt 200811.txt | sort | uniq -c | head -n 5
- 633 A
 - 1 AA
 - AARP
 - 1 ABBY
 - 41 ABC

- tr -sc 'A-Za-z' '\n' < nyt 200811.txt | sort | uniq -c | head
- **head** gives you the first 10 lines
- tail does the same with the end of the input
- (You can omit the "-n" but it's discouraged.)

Ex 2: Extended Counting Exercises

- 1. Merge upper and lower case by downcasing everything
 - Hint: Put in a second tr command 0

- 2. How common are different sequences of vowels (e.g., the sequences "ieu" or just "e" in "lieutenant")?
 - Hint: Put in a second tr command 0

Solutions

Merge upper and lower case by downcasing everything tr -sc 'A-Za-z' '\n' < nyt 200811.txt | tr 'A-Z' 'a-z' | sort | uniq -c or tr -sc 'A-Za-z' '\n' < nyt 200811.txt | tr

- '[:upper:]' '[:lower:]' | sort | uniq -c
- tokenize by replacing the complement of letters with newlines 1.
- 2. replace all uppercase with lowercase
- 3. sort alphabetically
- 4. merge duplicates and show counts

Solutions

How common are different sequences of vowels (e.g., ieu)

tr 'A-Z' 'a-z' < nyt_200811.txt | tr -sc 'aeiou' '\n' | sort | uniq -c

wels (e.g., xt | tr q -c

Sorting and reversing lines of text

- sort
- sort -f **Ignore case**
- sort -n Numeric order
- sort -r Reverse sort
- sort -nr Reverse numeric sort

• echo "Hello" | rev

Ex 3: Counting and sorting exercises

Find the 50 most common words in the NYT

• Hint: Use sort a second time, then head

Find the words in the NYT that end in "zz"

- Hint: Look at the end of a list of reversed words
- tr 'A-Z' 'a-z' < filename | tr -sc 'a-z' '\n' | rev | sort | rev | uniq –c

Ex 3 Counting and sorting solutions

Find the 50 most common words in the NYT

tr -sc 'A-Za-z' '\n' < nyt 200811.txt | sort | uniq -c | sort -nr | head -n 50

Find the words in the NYT that end in "zz"

tr -sc 'A-Za-z' '\n' < nyt 200811.txt | tr 'A-Z' 'a-z' | rev | sort | uniq -c | rev | tail -n 10

grep

Grep finds patterns specified as regular expressions grep rebuilt nyt 200811.txt

Conn and Johnson, has been rebuilt, among the first of the 222 move into their rebuilt home, sleeping under the same roof for the the part of town that was wiped away and is being rebuilt. That is to laser trace what was there and rebuilt it with accuracy," she home - is expected to be rebuilt by spring. Braasch promises that

grep

Grep finds patterns specified as regular expressions • globally search for regular expression and print

(A much easier way to find words ending in a pattern than the rev method we used before) Finding words ending in –ing: grep 'ing\$' nyt.words |sort | uniq -c

grep

grep is a filter – you keep only some lines of the input keep lines containing "gh" grep gh keep lines beginning with "con" grep '^con' grep 'ing\$' keep lines ending with "ing" keep lines NOT containing "gh" grep -v gh

grep versus egrep (grep –E)

egrep or grep -E [extended syntax] In egrep, +, ?, |, (, and) are automatically metacharacters In grep, you have to backslash them To find words ALL IN UPPERCASE: egrep '^[A-Z]+\$' nyt.words |sort|uniq -c ==grep '^[A-Z] \+\$' nyt.words |sort|uniq -c

(confusingly on some systems grep acts like egrep

Ex 4: Exercises on grep & wc

How many all uppercase words are there in this NYT file? How many 4-letter words?

How many different words are there with no vowels What subtypes do they belong to? 0

Type/instance distinction: different words (types) vs. instances (sometimes called "type/token" distinction but we now save "token" for BPE tokens)

Ex 4 Solutions on grep & wc

How many all uppercase words are there in this NYT file? grep -E '^[A-Z]+\$' nyt.words | wc How many 4-letter words? grep -E $\left[a-zA-Z\right]\left\{4\right\}$ nyt.words | wc How many different words are there with no vowels grep -v '[AEIOUaeiou]' nyt.words | sort | uniq | wc

Type/instance distinction: different words (types) vs. instances



Piping commands together can be simple yet powerful in Unix

It gives flexibility.

Traditional Unix philosophy: small tools that can be composed

Ex 5: Data Ethics

The text for today's lab comes from the New York Times Annotated **Corpus**, released by the Linguistic Data Consortium.

Understanding a dataset's origins gives us insight into its scope and limitations, which can affect the conclusions we draw from it.

Take a few minutes with your group to explore the origins of this text.

What types of texts, from what time period, genre, and language, are included in this dataset?

What types are not? Then, answer questions on the following slide

Data Ethics Questions

- 1. Imagine you are a researcher creating an educational platform as part of a high-school history course, and you trained your model solely on text from this corpus. Though the NYT is a respected publication, its content reflects its very particular context. What biases could arise from relying solely on this corpus for training your model? (For example, whose perspectives are predominantly included, and whose might be underrepresented or excluded?)
- 2. With your group, brainstorm ways to reduce the biases identified in part (1) when building this tool. How would you evaluate the effectiveness of these methods in reducing biases?

Bigrams = word pairs and their counts

Algorithm:

- 1. Tokenize by word
- 2. Create two almost-duplicate files of words, off by one line, using tail
- paste them together so as to 3. get word; and word; +1 on the same line

4. Count

Bigrams

- tr -sc 'A-Za-z' '\n' < nyt 200811.txt > nyt.words
- tail -n +2 nyt.words > nyt.nextwords
- paste nyt.words nyt.nextwords > nyt.bigrams
- head -n 5 nyt.bigrams

KBR said said Friday Friday the the global global economic



Ex 6: Bigrams

Find the 10 most common bigrams

• (Just for you to notice:) What part-of-speech pattern are most of them?

Ex 6 Solutions

Find the 10 most common bigrams

tr 'A-Z' 'a-z' < nyt.bigrams | sort | uniq -c | sort -nr | head -n 10

sed

sed is used to change strings in a file (larger changes) than 'tr')

- sed is line-based. You specify a regex substitution to make on each line
- (Optional: you can specify a subset of lines, by regex or line numbers)

For example, to change all cases of "George" to "Jane":

sed 's/George/Jane/' nyt 200811.txt less

Ex 7: sed exercises

1. Count frequency of word initial consonant sequences

- Take tokenized words
- Delete the first vowel through the end of the word 0
- Sort and count

2. Optional: Count word final consonant sequences

Ex 7 sed exercises solutions

- Count frequency of word initial consonant sequences tr "[:upper:]" "[:lower:]" < nyt.words | sed</pre> 's/[aeiou].*\$//' | sort | uniq -c
- Optional: Count word final consonant sequences
- tr "[:upper:]" "[:lower:]" < nyt.words | sed</pre> 's/^.*[aeiou]//' | sort | uniq -c | sort -rn | less

Extra Credit – Secret Message

- Now, let's get some more practice with Unix!
- The answers to the extra credit exercises will reveal a secret • message.
- We will be working with the following text file for these exercises:

https://web.stanford.edu/class/cs124/lec/secret ec.txt

- To receive credit, enter the secret message here:
- https://docs.google.com/forms/d/e/1FAIpQLSdrr1p31AGkgZJSu
- 34 7-sLtxbobJIsPdkq5KXe0n7jKZosGWNAw/viewform

Fxtra Credit Exercise 1

Find the 2 most common words in secret ec.txt containing the letter e.

Your answer will correspond to the first two words of the secret message.

Extra Credit Exercise 2

Find the 2 most common bigrams in secret ec.txt where the second word in the bigram ends with a consonant.

Your answer will correspond to the next four words of the secret message.

Fxtra Credit Exercise 3

Find all 5-letter-long words that only appear once in secret ec.txt.

Concatenate (by hand) your result. This will be the final word of the secret message.