BINF2111 - Introduction to Bioinformatics Computing UNIX 101 - enter the coding zone



Lecture 5 - Tuesday Sep 7th, 2021

RAW Lab

Learning Objectives

- Review bonus

- Comments about data and formats

- Sort/uniq/cut commands
- PATHS
- Quiz 5

- Delete all the empty lines in the empty lines file with

- Delete all the 'all white space' with grep

- Delete all the empty lines in the empty lines file with
- → grep: grep -v -e '^\$' file

- Delete all the 'all white space' with grep

- Delete all the empty lines in the empty lines file with
- → grep: grep -v -e '^\$' file
- \rightarrow awk: awk '!/^\$/' file
- Delete all the 'all white space' with grep

- Delete all the empty lines in the empty lines file with
- → grep: grep -v -e '^\$' file
- → awk: awk '!/^\$/' file
- Delete all the 'all white space' with grep
- → grep: grep -v -e '^[[:space:]]*\$' file

- Delete all the empty lines in the empty lines file with
- → grep: grep -v -e '^\$' file
- → awk: awk '!/^\$/' file
- Delete all the 'all white space' with grep
- → grep: grep -v -e '^[[:space:]]*\$' file
- → awk: awk 'NF > 0' file
- → sed 's/\t/,/g' empty_lines.txt >empty_lines.csv

Also, in python (think Pandas)

Any one try in Python Pandas.

Bonus 4 – Python Pandas

- First use sed to convert from tsv (tab delim) to csv (column delim)

import pandas as pd

```
df = pd.read_csv('data.csv')
new_df = df.dropna() or df.dropna(inplace = True)
print(new_df.to_string()) or print(df.to_string())
```

General comments - data

UNDERSTAND YOUR DATA

- Expectations Format, Types, Domains, Uniqueness, etc.
- Assumptions Does it make sense to use data from X in a new context Y? Restrictions Licenses, Embargos, etc.

STANDARDIZE YOUR DATA

- Values ID mapping, Unit conversion, Adding/Removing Prefixes/Suffixes, etc.
- Columns Adding, Removing, Rearranging, Merging, Splitting, etc.
- Rows Filtering data (duplicates, missing data, not relevant to task, etc)

RECORD YOUR ANALYSIS PROCESS

- Your computer is a lab Keep a lab notebook!
- Use Version Control like GitHub or BitBucket.

General comments - data formats

Column-oriented data

- Spreadsheets, CSV, Tabular, delimited
- Fixed position
- One line per record, fixed set of fields (columns)

Key-Value data

- Multiple lines per record
- Variable number of fields/values

Hierarchical data (XML, JSON, Ontologies, etc)

- Nested usually requires a more complex parser
- Usually follow a well-defined schema.

Web-based API resources

- Issue specific commands to get different types of data
- **Infinite proprietary formats**
- Usually harder to parse but generally act like other types

General comments: Column-Oriented

Easy to browse/explore

- MOST files are this type

Best for databases

- Delimiters are tabs or commas between fields
- 1 line = 1 record
- Each record has fixed set of fields

Complex code not needed:

- Excel can handle small datasets(<1 million rows)
- UNIX can handle the rest

CSV

david,abdul,xi,bill mary,david,bill,abdul wang,abdul,xi,david

TSV

Mary david bill abdul Wang abdul xi david

David abdul xi bill

General comments - XML formats

Machine-readable

- Many tools and libraries available to parse it

Can represent complex structures

- multiple values
- nested structures

Not trivial for non-coders

```
<?xml version="1.0"?>
<!DOCTYPE Entrezgene-Set PUBLIC "-//NLM//DTD NCBI-Entrezgene, 21st Ja
<Entrezgene-Set>
<Entrezgene>
  <Entrezgene track-info>
    <Gene-track>
      <Gene-track geneid>4336</Gene-track geneid>
      <Gene-track status value="live">0</Gene-track status>
      <Gene-track_create-date>
        <Date>
          <Date std>
            <Date-std>
              <Date-std_year>1998</Date-std_year>
              <Date-std_month>8</Date-std_month>
              <Date-std day>27</Date-std day>
            </Date-std>
          </Date std>
        </Date>
      </Gene-track_create-date>
      <Gene-track_update-date>
        <Date>
          <Date std>
            <Date-std>
              <Date-std_year>2016</Date-std_year>
              <Date-std_month>12</Date-std_month>
              <Date-std_day>6</Date-std_day>
            </Date-std>
          </Date std>
        </Date>
      </Gene-track update-date>
    </Gene-track>
  </Entrezgene track-info>
  <Entrezgene_type value="protein-coding">6</Entrezgene_type>
  <Entrezgene source>
    <BioSource>
      <BioSource genome value="genomic">1</BioSource genome>
      <BioSource origin value="natural">1</BioSource origin>
      <BioSource org>
        <0rg-ref>
          <Org-ref_taxname>Homo sapiens/Org-ref_taxname>
          <0rg-ref_common>human</0rg-ref_common>
          <0rg-ref db>
            <Dbtag>
              <Dbtag_db>taxon</Dbtag_db>
              <Dbtag tag>
                <Object-id>
                  <Object-id_id>9606</Object-id_id>
                </Object-id>
              </Dbtag tag>
            </Dbtag>
          </0rg-ref db>
          <0rg-ref_syn>
            <Org-ref_syn_E>humans</Org-ref_syn_E>
            <Org-ref_syn_E>man</Org-ref_syn_E>
          </0rq-ref_syn>
          <Org-ref_orgname>
            <OrgName>
              <OrgName_name>
```

General comments: Key-value formats

Easy to read

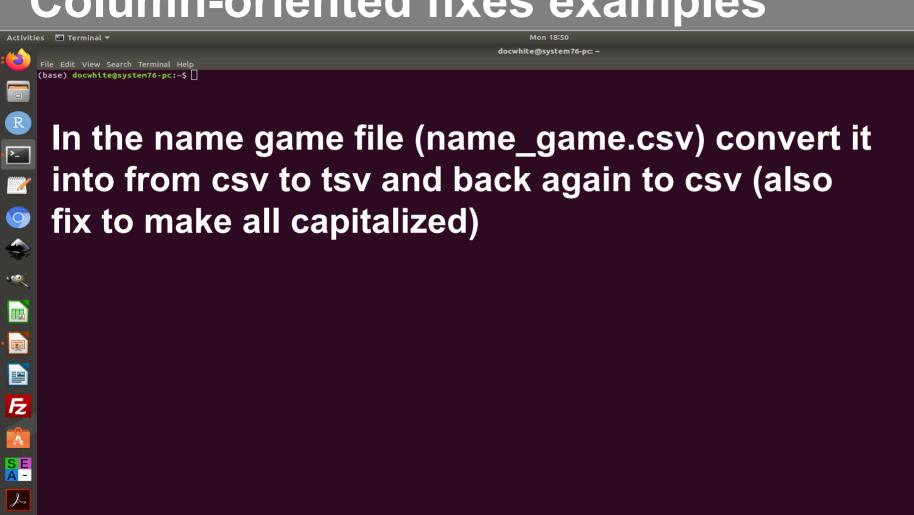
Straightforward to parse

- Delimiter for each Record
- One key-value pair per line
- Supports multiple values

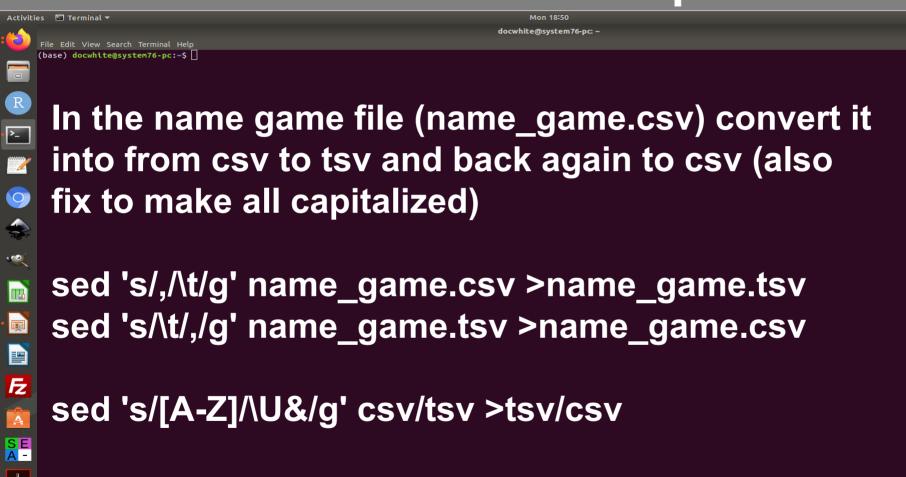
Novice coders can handle it!

```
remark: cvs version: use data-version
ontology: go
[Term]
id: GO:0000001
name: mitochondrion inheritance
namespace: biological process
def: "The distribution of mitochondria, including the mitocho
PMID:10873824, PMID:113897641
exact synonym: "mitochondrial inheritance" []
is a: GO:0048308 ! organelle inheritance
is a: GO:0048311 ! mitochondrion distribution
[Term]
id: GO:0000002
name: mitochondrial genome maintenance
namespace: biological process
def: "The maintenance of the structure and integrity of the m
is a: GO:0007005 ! mitochondrion organization
[Term]
id: GO:0000003
name: reproduction
namespace: biological process
alt id: GO:0019952
alt id: GO:0050876
def: "The production of new individuals that contain some por
subset: goslim chembl
subset: goslim generic
subset: goslim pir
subset: goslim plant
subset: gosubset prok
exact synonym: "reproductive physiological process" []
xref analog: Wikipedia: Reproduction
is a: GO:0008150 ! biological process
[Term]
id: GO:0000005
name: obsolete ribosomal chaperone activity
namespace: molecular function
def: "OBSOLETE. Assists in the correct assembly of ribosomes
PMID:121509131
comment: This term was made obsolete because it refers to a c
exact synonym: "ribosomal chaperone activity" []
is obsolete: true
consider: GO:0042254
consider: GO:0044183
consider: GO:0051082
```

Column-oriented fixes examples



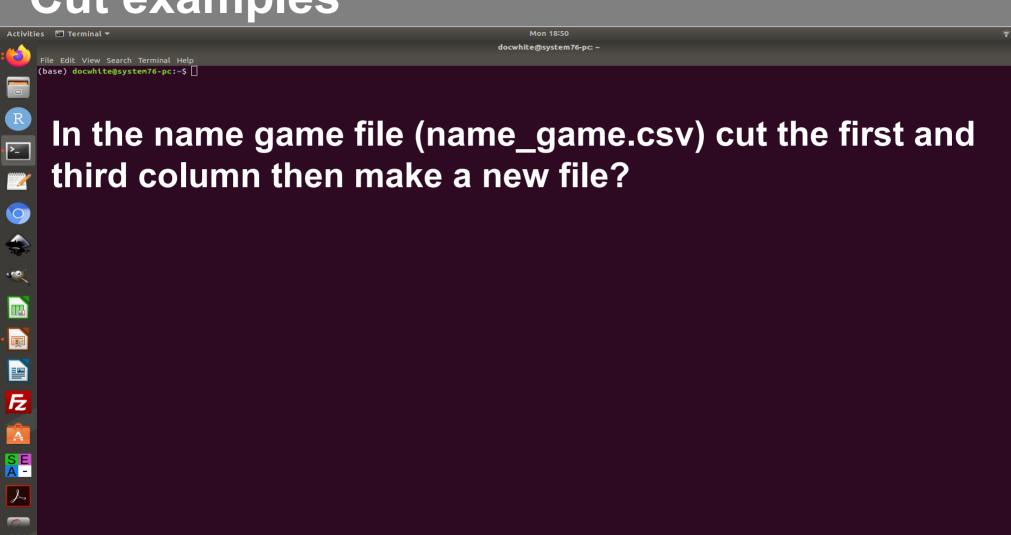
Column-oriented fixes examples

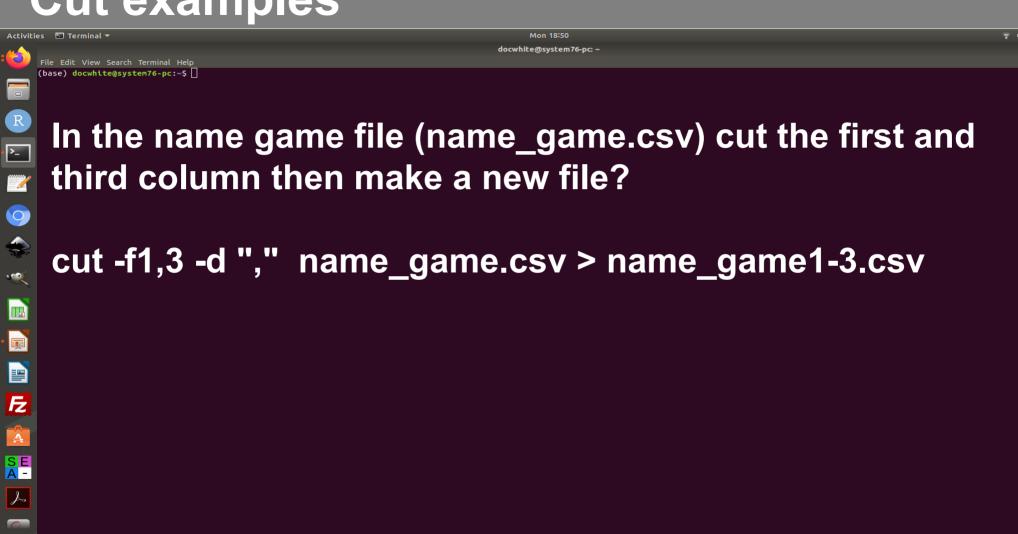


cut – syntax anatomy UNIX's scissors

cut [options] file.txt

- -d (--delimiter) "," set field delimiter (default tab)
- -f (--fields=LIST) Select by specifying a field
 - -f 2 select a field to cut (left is 1)
 - -f 2-8,12 select multiple fields to cut
- -b (--bytes=LIST) Select by specifying a byte
- -c (--characters=LIST) Select by specifying a character
- --complement Complement the selection.
- -s (--only-delimited) suppress non-matches









sort – syntax anatomy of sort

sort [options] file.txt

- -t "t" set the delimiter when using -k default is non-blank to blank transition
- -k 3 sort column #3 (left is 1)
- -k 2,3 sort multiple columns
- -n sort numerically
- -r reverse sort order
- -u drop duplicates from the result
- -b, --ignore-leading-blanks, ignore leading blanks
- -d, --dictionary-order consider only blanks and alphanumeric characters
- -f, --ignore-case fold lower case to upper case characters
- No options: sort alphabetically from leftmost character.

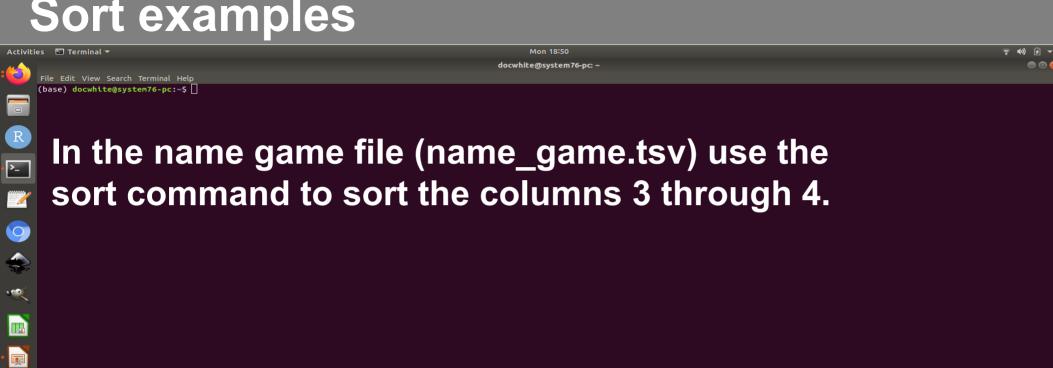
Æ



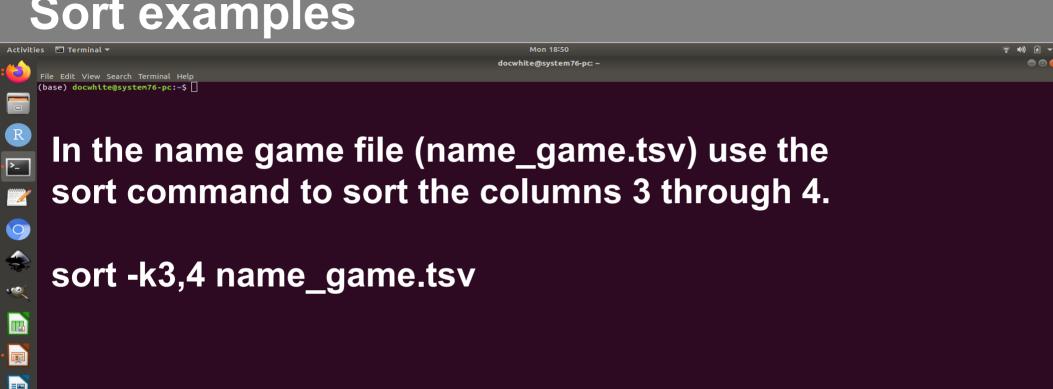


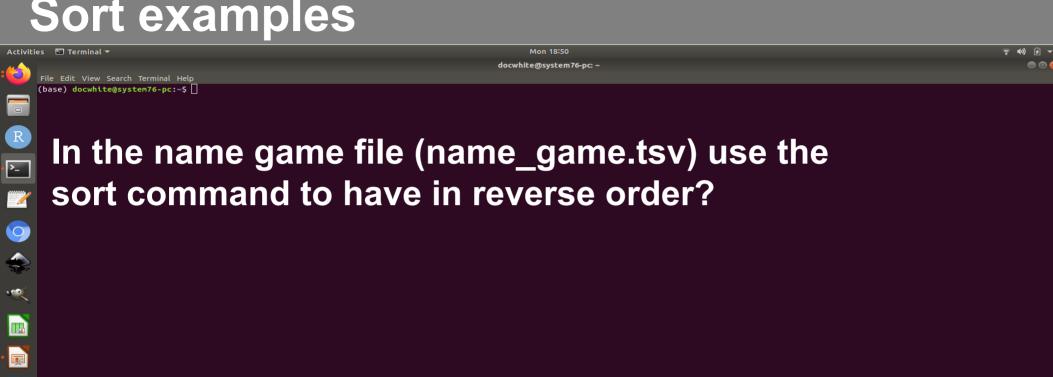
As tab (t) is default

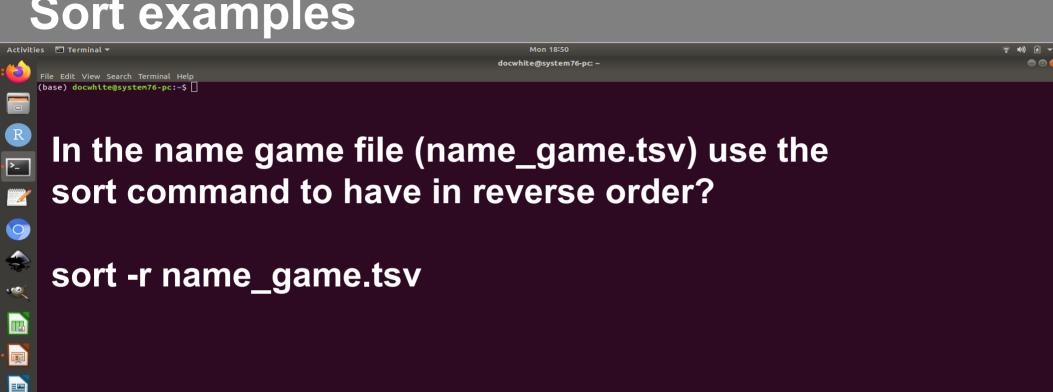
Æ

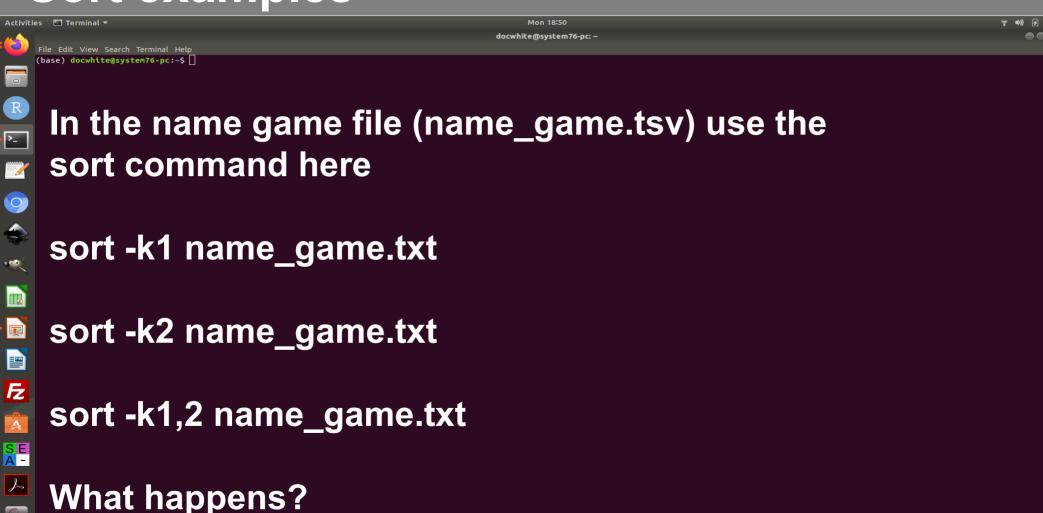


Æ





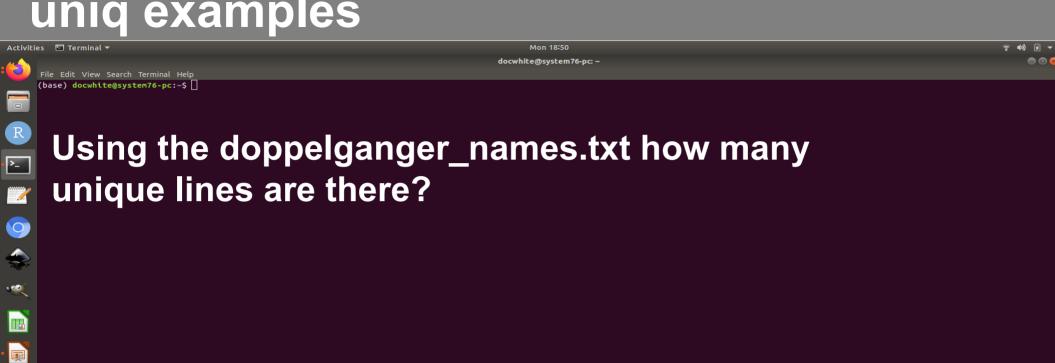


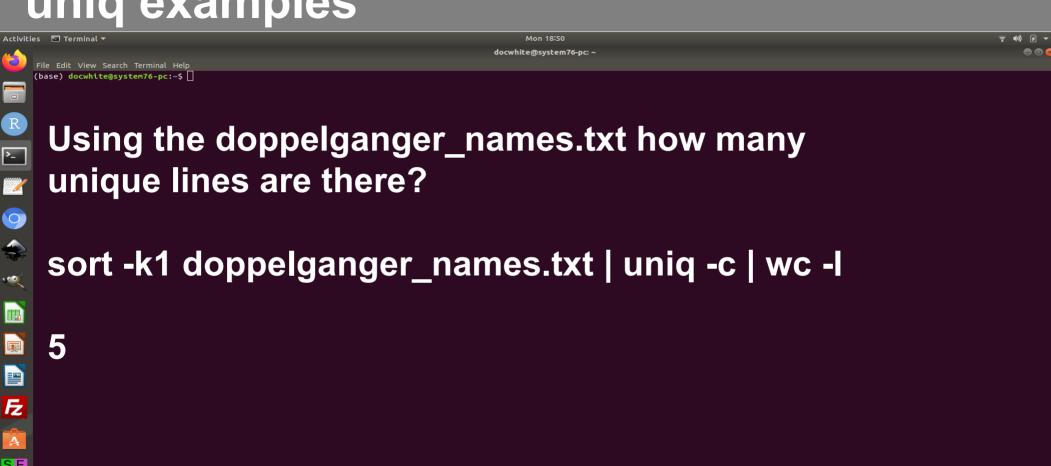


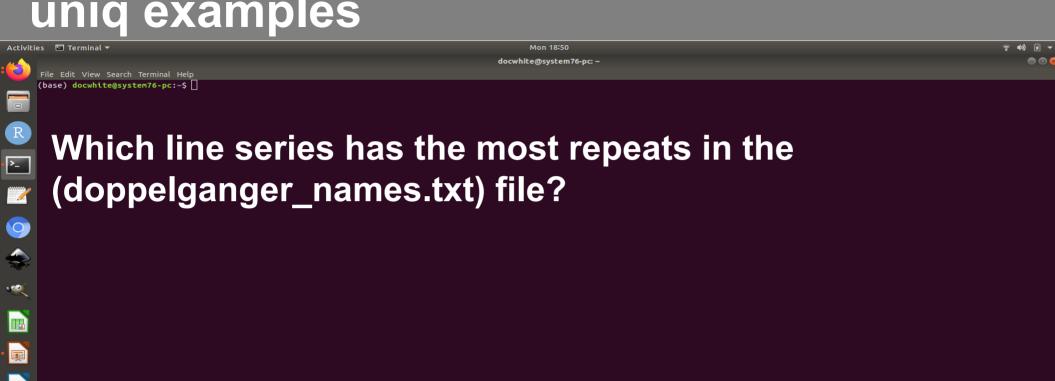
uniq – syntax anatomy of uniq

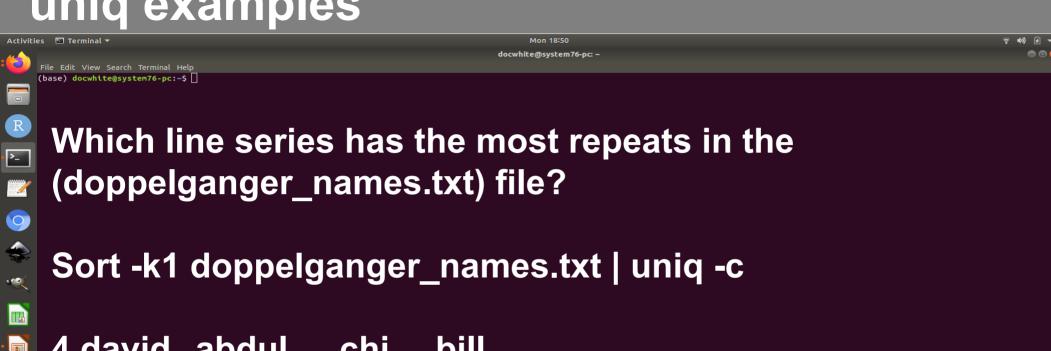
uniq [options] file.txt

- Only works on sorted files and adjacent lines!
- -c count lines for each unique value
- -d only report duplicated lines
- -u only report non-duplicated lines
- No options: drop all duplicated lines (keeps 1 copy)







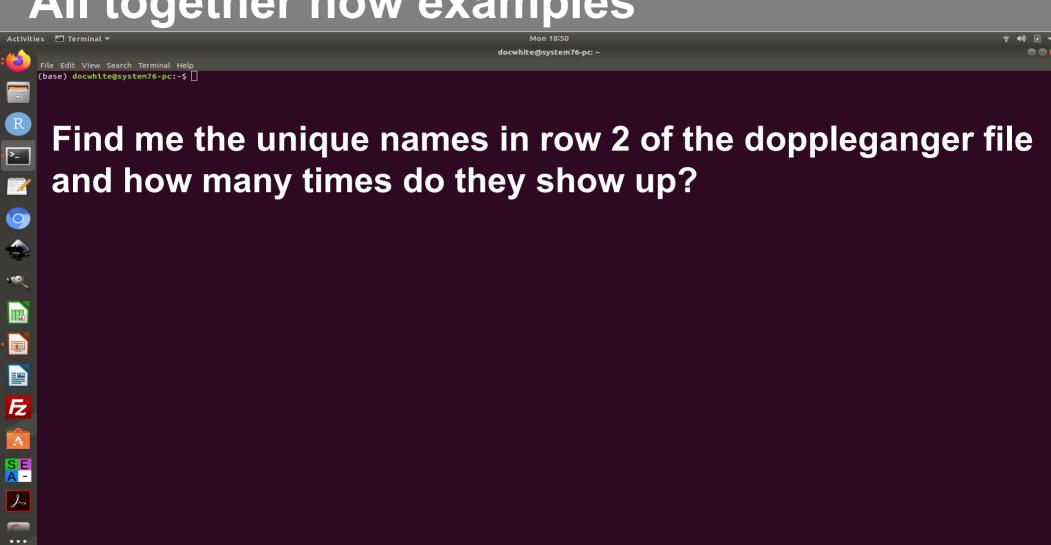


- 4 david abdul chi bill bill abdul david 4 mary
- SE A-

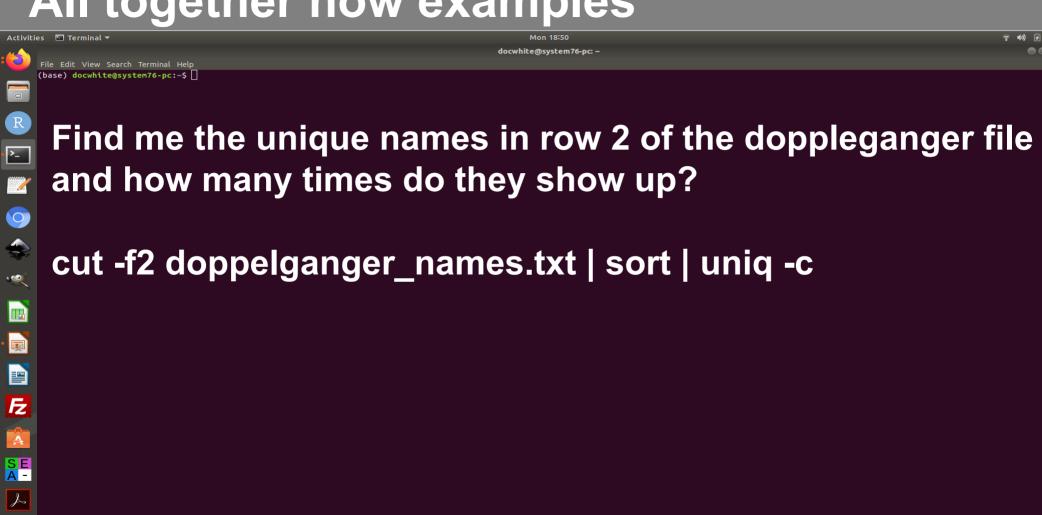




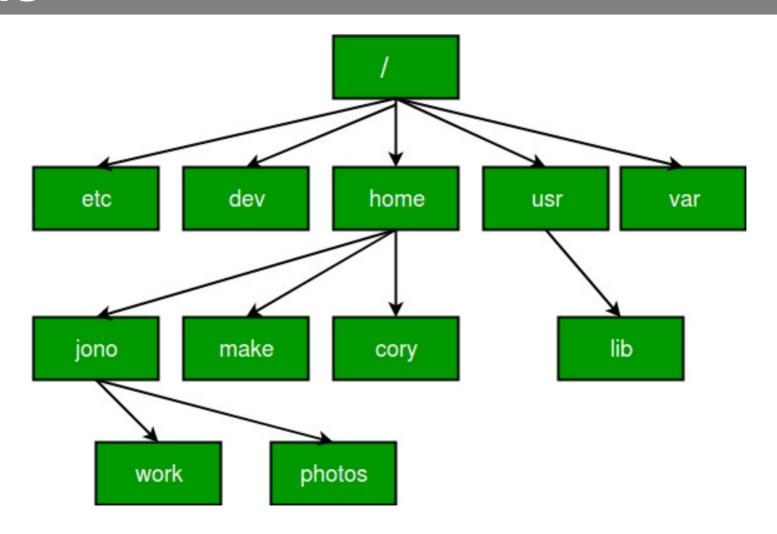
All together now examples



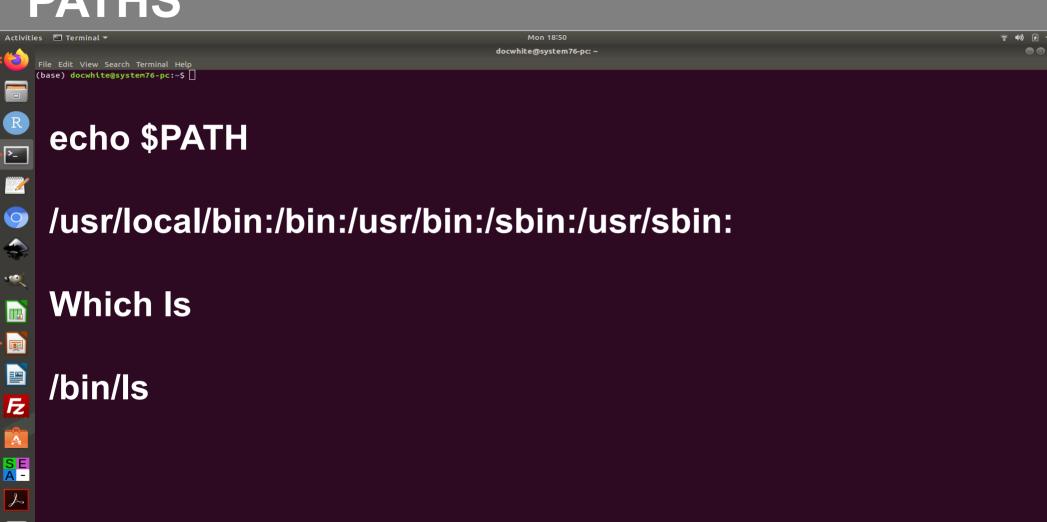
All together now examples



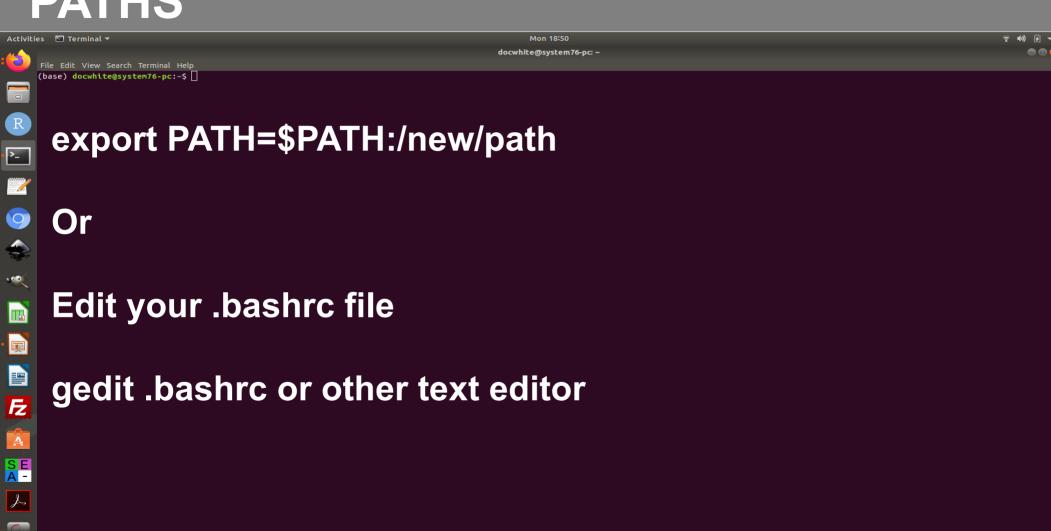
PATHS



PATHS



PATHS



- Convert name game file (name_game.csv) to tsv with:

 \rightarrow tr

 \rightarrow awk

Quiz 5

- On canvas now