# Data Analysis for Molecular Biology and Biochemistry

#### **MBB 110**



Introductory data analysis focusing on molecular biology data sets and examples and including basic programming skills using Python and basic statistics skills using R. Prerequisite: MATH 12 or equivalent is recommended. Students with credit for MBB 243 may not take this course for further credit. CMPT 120 will be accepted in lieu of MBB 110.

## **Topics**

- Flavours of data in molecular biology and biochemistry
- Genomic data
- Fundamentals of R and Python
- Regular expressions and patterns
- Quantitative DNA/RNA sequence analysis
- Exploratory data analysis in the Tidyverse
- Generic visualization methods using ggplot2
- · Advanced visualizations for molecular biology

#### **INSTRUCTOR:**

#### Sophie Sneddon



15000

sort sorts its input

\$ sort names.txt

the default sort is alphabetical.

sort-n		
numeric sort		
'sort' order	'sort -n'	order
12	12	
15000	48	$\mathbb{D}$
,, 48	96	IJ
6020	6020	

96

Sort -h: human sort

'sort -n' order | 'sort -h' order

156 | 45 K

130 M | 30 M | 156

145 K | 156 | 156

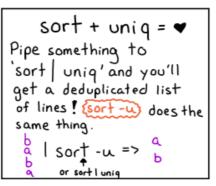
2006 | 2006

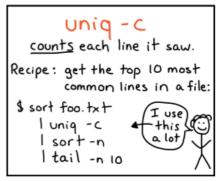
Useful example:

du-sh \* | sort -h

# Uniq removes duplicates

a notice there
b are still 2
b 'a's! uniq
a only uniquifies
c adjacent
matching lines





## **Group Data**

dplyr::group\_by(iris, Species)

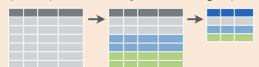
Group data into rows with the same value of Species.

dplyr::ungroup(iris)

Remove grouping information from data frame.

iris %>% group\_by(Species) %>% summarise(...)

Compute separate summary row for each group.



iris %>% group\_by(Species) %>% mutate(...)
Compute new variables by group.

