

R Programming For Natural Resource Professionals

Lecture 5 Data Wrangling II: Joining and advanced dplyr

Polishing R markdowns

Updated criteria for R markdown documents

- 1) Display resulting tibbles/tables using kableExtra
`tib %>% kbl() %>% kable_styling()`
[Consult excellent kableExtra vignette for options](#)
- 2) Final document emphasizes code and result. No messages, errors, etc.
- 3) Goal is to generate publication quality documents
- 4) Get creative and have fun with it.

Working with row names

A	B	C
1	a	t
2	b	u
3	c	v

→

A	B
1	a
2	b
3	c

`tibble::column_to_rownames()`

A	B
1	a
2	b
3	c

→

C	A	B
1	a	t
2	b	u
3	c	v

`tibble::rownames_to_column()`

Adding new data



`dplyr::mutate`

Add a new variable or change an existing variable

```
mutate(newVar = [calculation])
```

Adding new data



`dplyr::transmute`

Compute a new variable while dropping the others

```
transmute(newVar = [calculation])
```

Rename variables



`dp1yr::rename`

Change the name of a variable

```
rename(newName = oldName)
```

Combining data sets

$+$

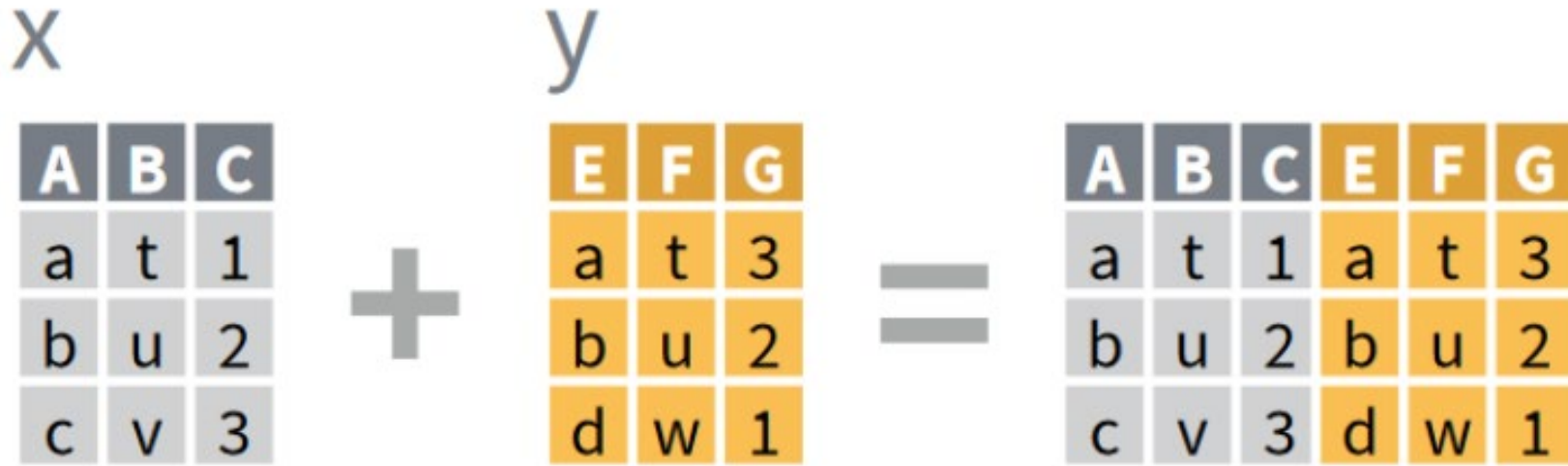
x	<table><tr><th>A</th><th>B</th><th>C</th></tr><tr><td>a</td><td>t</td><td>1</td></tr><tr><td>b</td><td>u</td><td>2</td></tr></table>	A	B	C	a	t	1	b	u	2	<table><tr><th>A</th><th>B</th><th>C</th></tr><tr><td>c</td><td>v</td><td>3</td></tr><tr><td>d</td><td>w</td><td>4</td></tr></table>	A	B	C	c	v	3	d	w	4
A	B	C																		
a	t	1																		
b	u	2																		
A	B	C																		
c	v	3																		
d	w	4																		
y																				

A	B	C
a	t	1
b	u	2
c	v	3
d	w	4

`bind_rows()`

Returns one tibble pasted above the other

Combining data sets



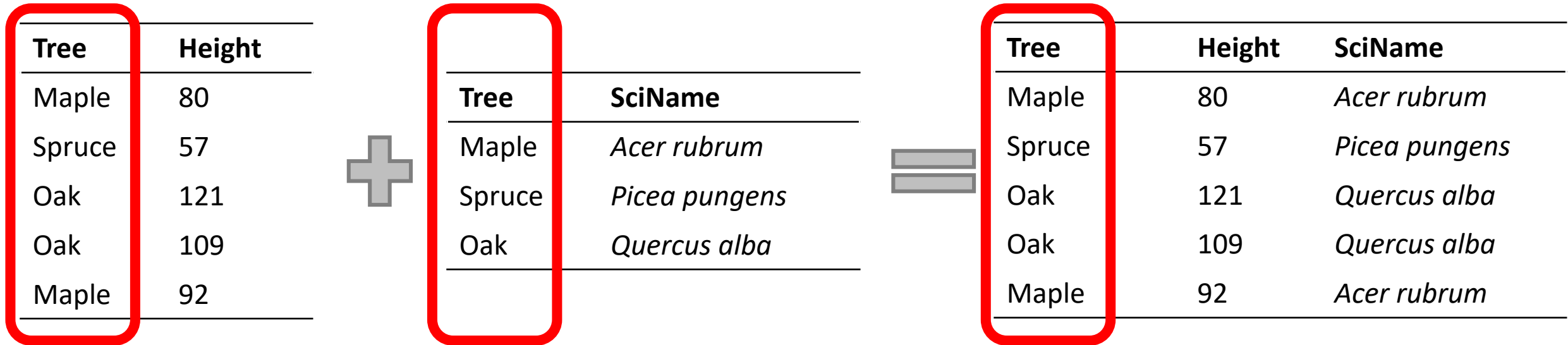
`bind_cols()`

Returns one tibble pasted next to the other.

Important: Not for combining tibbles with the same variables!

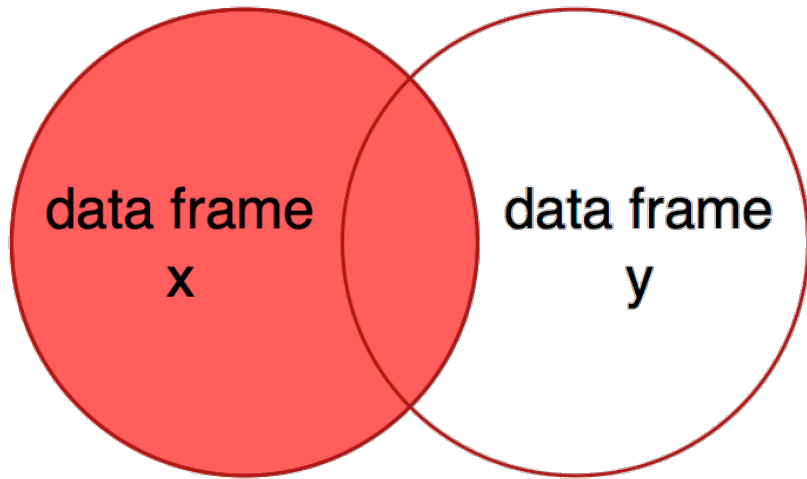
Joining data sets

- Merge data based on given criteria
 - “Relational merging”



Combining data sets

`left_join()`



Data frame x

ID	X1
1	a1
2	a2

Data frame y

ID	X2
2	b1
3	b2

Result

ID	X1	X2
1	a1	NA
2	a2	b1

- Key variable in example is “ID”
- Returns all observations from x and all variables from x and y.
- Observations in x with no match in y will be populated with NAs

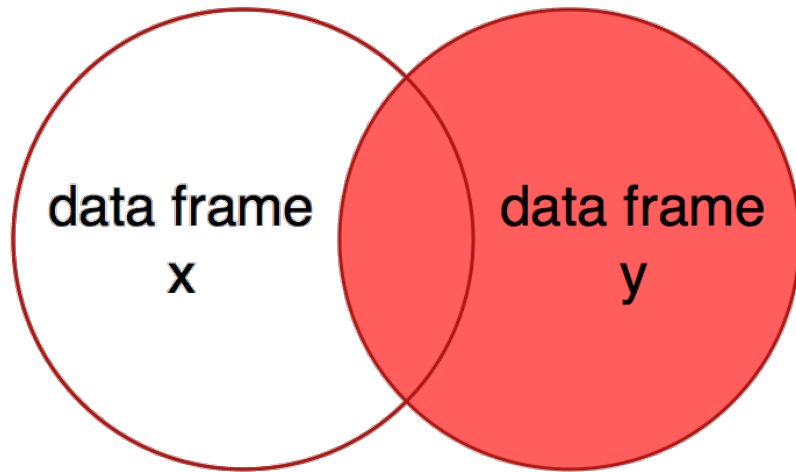
Combining data sets

`left_join(x, y)`

1	x1	1	y1
2	x2	2	y2
3	x3	4	y4

Combining data sets

`right_join()`



Data frame x	
ID	X1
1	a1
2	a2

Data frame y	
ID	X2
2	b1
3	b2

Result		
ID	X1	X2
2	a2	b1
3	NA	b2

- Key variable in example is “ID”
- Returns all observations from y and all variables from x and y.
- Observations in y with no match in x will be populated with NAs

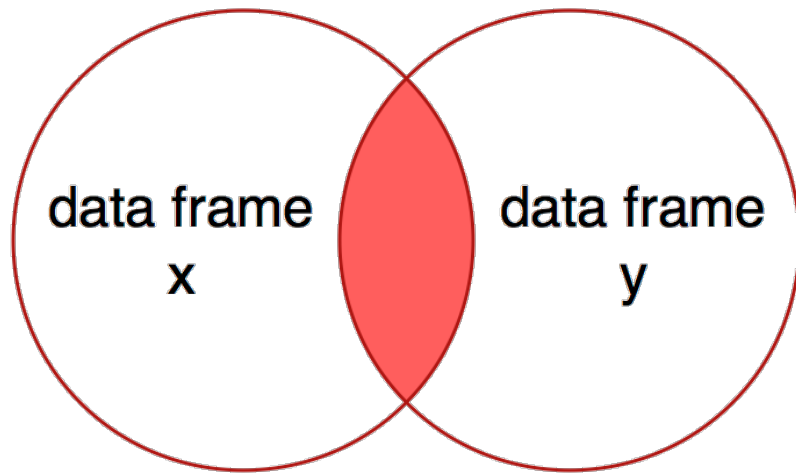
Combining data sets

`right_join(x, y)`

1	x1	1	y1
2	x2	2	y2
3	x3	4	y4

Combining data sets

`inner_join()`



Data frame x

ID	X1
1	a1
2	a2

Data frame y

ID	X2
2	b1
3	b2

Result

ID	X1	X2
2	a2	b1

Key variable in example is “ID”

Returns all observations from x with matching observations in y

Returns all columns of x and y

If there are multiple matches between x and y, all combinations are returned

Combining data sets

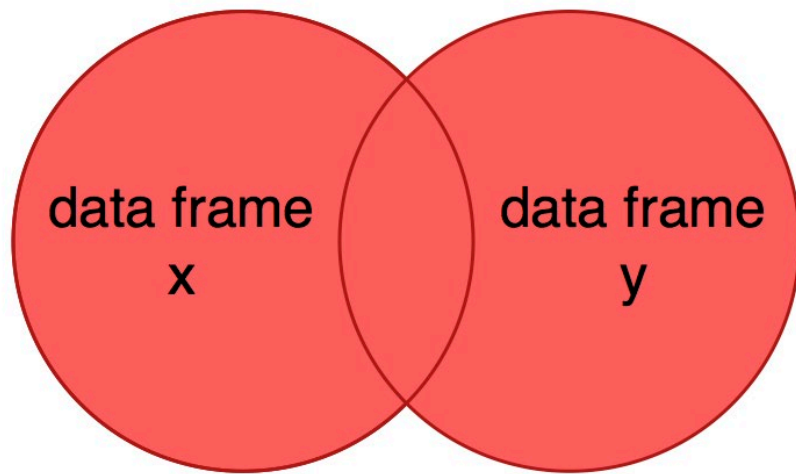
`inner_join(x, y)`

1	x1
2	x2
3	x3

1	y1
2	y2
4	y4

Combining data sets

`full_join()`



Data frame x	
ID	X1
1	a1
2	a2

Data frame y	
ID	X2
2	b1
3	b2

Result		
ID	X1	X2
1	a1	NA
2	a2	b1
3	NA	b2

Returns all observations and variables of both x and y
When not matching, returns NA

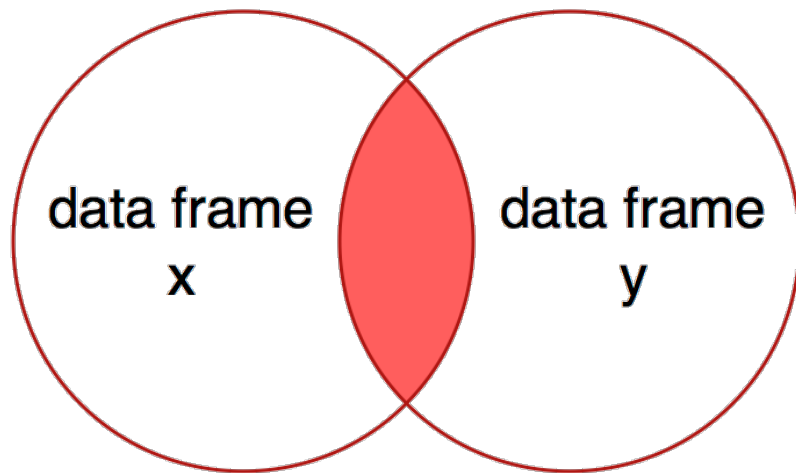
Combining data sets

`full_join(x, y)`

1	x1	1	y1
2	x2	2	y2
3	x3	4	y4

Combining data sets

`semi_join()`



Data frame x	
ID	X1
1	a1
2	a2

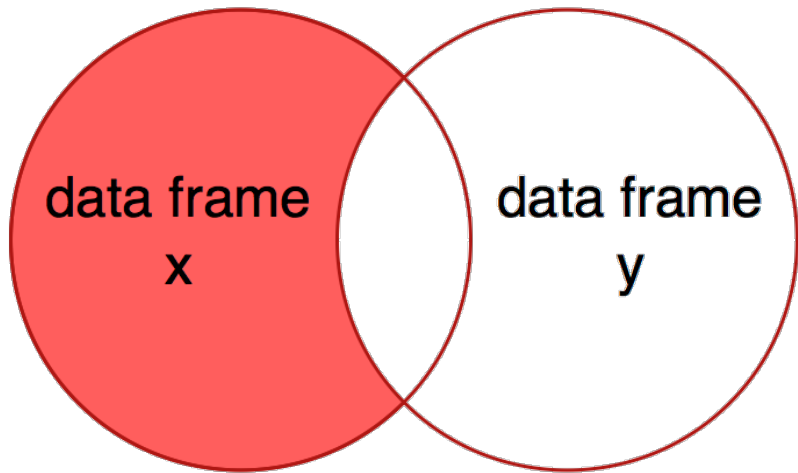
Data frame y	
ID	X2
2	b1
3	b2

Result	
ID	X1
2	a2

- A type of “filtering join”
- Returns observations from x that **do** also occur in y.
- Differs from `inner_join()` because it doesn't retain variables in y

Combining data sets

`anti_join()`



Data frame x	
ID	X1
1	a1
2	a2

Data frame y	
ID	X2
2	b1
3	b2

Result	
ID	X1
1	a1

- A type of “filtering join”
- Returns observations from x that **do not** also occur in y.

Combine Data Sets

a				b			
x1	x2			x1	x3		
A	1			A	T		
B	2			B	F		
C	3			D	T		

+

=

Combining data sets

- What if the key variable isn't the same in x and y?

```
tib1 %>%
```

```
  left_join(tib2, by = c("fish" = "species"))
```

Tib1	
fish	length
Brook trout	12
Walleye	15
Walleye	16

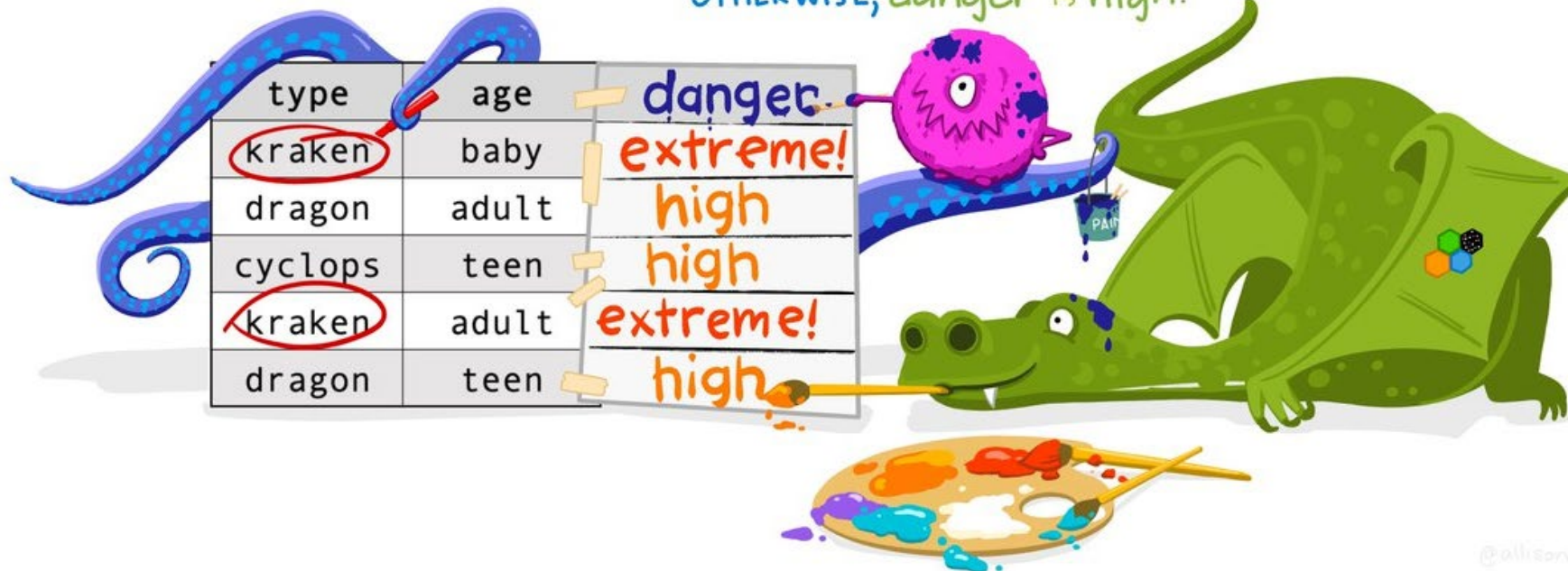
Tib2	
species	genus_species
Brook trout	<i>Salvelinus fontinalis</i>
Walleye	<i>Sander vitreus</i>

Create new variables using conditional statements

dplyr::case_when() IF ELSE...
(but you love it?)

df %>% ^{ADD COLUMN 'danger'}
mutate(danger = case_when(type == "kraken" ^{IF type is kraken} ^{THEN} "extreme!",
T ~ "high")) ^{danger is extreme!}

OTHERWISE, danger is high.



Paired programming

Exercise 1

Step 1: Use a join operation to generate a tibble containing all the survey data, completed with the data in the species tibble.

Step 2: Combine the genus and species column into a single column separated by an underscore. Retain the original columns.

Step 3: Subset your data to include only observation of *Dipodomys merriami*.

Step 4: Calculate the mean and standard deviation of the hindfoot length of this species for each year of the data.

Step 5: Within the pipeline, filter the data set to just the maximum mean hindfoot length

Paired programming

Exercise 2

Not all survey observations have a corresponding value in the species table. Use a filtering join to determine which do. Then determine how many different values of `species_id` occur.