

# Using regular expressions

Lauren Ponisio

## Conservation/ecology Topics

- Species distributions

## Computational Topics

- Use regular expressions to clean and categorize data
- 

## Part 1: Oregon bee atlas data exploration

Import the OBA data using your favorite parsing function, name the data oba.

a.

```
oba <- read.csv("Data/OBA_2018-2023.csv")
```

- b. Examine the unique entries of 'Associated.plant' using any function you find useful. What are at least two patterns in the associated taxa string what should be removed if we want consistent plant names? (Make a list together as a class). Only print the first 10 here to avoid having a giant output.

```
head(oba)
```

```
##              Observation.No. Voucher.No. user_id    user_login
## 1 Andony_Melathopoulos:18.001.001      429964 amelathopoulos
## 2 Andony_Melathopoulos:18.002.001      429964 amelathopoulos
## 3 Andony_Melathopoulos:18.002.002      429964 amelathopoulos
## 4 Andony_Melathopoulos:18.002.003      429964 amelathopoulos
## 5 Andony_Melathopoulos:18.002.004      429964 amelathopoulos
## 6 Andony_Melathopoulos:18.002.005      429964 amelathopoulos
## Collector...First.Name Collector...First.Initial Collector...Last.Name
## 1           Andony           A.           Melathopoulos
## 2           Andony           A.           Melathopoulos
## 3           Andony           A.           Melathopoulos
## 4           Andony           A.           Melathopoulos
## 5           Andony           A.           Melathopoulos
## 6           Andony           A.           Melathopoulos
```

```

##      Collectors taxon_kingdom_name Associated.plant...genus..species url
## 1 A.Melathopoulos
## 2 A.Melathopoulos
## 3 A.Melathopoulos
## 4 A.Melathopoulos
## 5 A.Melathopoulos
## 6 A.Melathopoulos
##   Sample.ID Specimen.ID Collection.Day.1 Month.1   MonthJul MonthAb Year.1
## 1           NA           18      iii    March        3   2018
## 2           NA           20      iii    March        3   2018
## 3           NA           20      iii    March        3   2018
## 4           NA           20      iii    March        3   2018
## 5           NA            2      ix September        9   2018
## 6           NA            2      ix September        9   2018
##   Collection.Date Time.1 Collection.Day.2 Month.2 Year.2 Collection.Day.2.Merge
## 1      3/18/2018
## 2      3/20/2018
## 3      3/20/2018
## 4      3/20/2018
## 5      9/2/2018
## 6      9/2/2018
##   Time.2      Collection.ID Position.of.1st.digit Collection.No. Sample.No.
## 1      A Melathopoulos                        1          1
## 2      A Melathopoulos                        2          1
## 3      A Melathopoulos                        2          2
## 4      A Melathopoulos                        2          3
## 5      A Melathopoulos                        2          4
## 6      A Melathopoulos                        2          5
##   Country State County                               Location
## 1      USA Oregon Benton                        Corvallis, NW Orchard Ave
## 2      USA Oregon Benton                        Corvallis, NW Orchard Ave
## 3      USA Oregon Benton                        Corvallis, NW Orchard Ave
## 4      USA Oregon Benton                        Corvallis, NW Orchard Ave
## 5      USA Oregon Clatsop Clatskanie, Big Creek Mainline, Knob Point Road
## 6      USA Oregon Clatsop Clatskanie, Big Creek Mainline, Knob Point Road
##   Abbreviated.Location Collection.Site.Description      Team
## 1      Astoria Maggie Johnson Rd                        Melathopoulos
## 2 Big Crk. Mainline Knob Pt Rd                        Melathopoulos
## 3 Big Crk. Mainline Knob Pt Rd                        Melathopoulos
## 4 Big Crk. Mainline Knob Pt Rd                        Melathopoulos
## 5 Big Crk. Mainline Knob Pt Rd                        Melathopoulos
## 6 Big Crk. Mainline Knob Pt Rd                        Melathopoulos
##   Habitat Elevation..m. Dec..Lat. Dec..Long. X Collectionmethod
## 1           44.556   -123.285 NA                      Net
## 2           44.567   -123.283 NA                      Net
## 3           44.567   -123.283 NA                      Net
## 4           44.567   -123.283 NA                      Net
## 5           46.102   -123.506 NA                      Net
## 6           46.102   -123.506 NA                      Net
##   Collection.method.merge.field Associated.plant...family
## 1
## 2
## 3
## 4

```

```

## 5
## 6
## Associated.plant...genus..species.1 Associated.plant...Inaturalist.URL
## 1
## 2
## 3
## 4
## 5
## 6
## Associated.plant Assoc.plant.merge.field Collectors.1
## 1 Andony Melathopoulos
## 2 Andony Melathopoulos
## 3 Andony Melathopoulos
## 4 Andony Melathopoulos
## 5 Andony Melathopoulos
## 6 Andony Melathopoulos
## Collector.1.abreviation Collector.2 Collector.3 Genus Species sex caste
## 1 A Melathopoulos NA NA
## 2 A Melathopoulos NA NA
## 3 A Melathopoulos NA NA
## 4 A Melathopoulos NA NA
## 5 A Melathopoulos NA NA
## 6 A Melathopoulos NA NA
## vol.det.Genus vol.det.Species vol.det.sex.caste Determined.By Date.Determined
## 1 NA
## 2 NA
## 3 NA
## 4 NA
## 5 NA
## 6 NA
## Verified.By Other.Determiner.s. Other.Dets.Sci..Name.s. Other.Dets..Date.s.
## 1 NA NA NA
## 2 NA NA NA
## 3 NA NA NA
## 4 NA NA NA
## 5 NA NA NA
## 6 NA NA NA
## Additional.Notes X.1
## 1 NA
## 2 NA
## 3 NA
## 4 NA
## 5 NA
## 6 NA

```

```

# examine unique entries
head(unique(oba$Associated.plant))

```

```

## [1] "" "Salix"
## [3] "Arctostaphylos densiflora" "Lithodora diffusa"
## [5] "Ceanothus gloriosus" "Ceanothus sp."

```

```
# what are two patterns in the associated taxa that should be removed?
# the ones that have common names after the genus/species name?, where sp doesn't have a period
```

1. Sometimes there is only one word, often family, genus, or common name. We will sort these out by creating a column for plant resolution.
2. The common name is sometimes listed after the scientific name in (), we can strip this out.
3. Sometimes there is a list of plant names, cannot do much with those but drop them.
4. Some are blank (no name), we will drop these.
5. One says “net”
6. Sometimes a genus is followed by an sp. and sometimes not. I assume this can be solved with a regular expression, ignoring the rows with a family name, but I could not work out how.
7. There are a few with “genus XX” or “genus XX”

There may be more I am missing. :/

In week in lecture last I used a brute force pattern to remove some of these issues so we could plot them as a network. Now that we are familiar with regular expressions we can do better.

- c. Work together as a class to resolve the issues you listed with the associated taxa column using any function combination that uses regular expressions. You can reassign the contents of the column `Associated.plant` or create a new column. Return the sorted, unique values, ex: `sort(unique(oba$Associated.plant))`. Leave the plants resolved only to genus or family for later.

I have removed a really strange issue with special characters (R converted an apostrophe into a special character) to start things off.

```
## Remove the special character
oba$Associated.plant <- str_replace_all(oba$Associated.plant, "\\x92", "")

## To check that it worked
sort(unique(oba$Associated.plant))[1:10]
```

```
## [1] "" "Abelia sp."
## [3] "Abronia latifolia" "Acer circinatum"
## [5] "Acer macrophyllum" "Acer palmatum"
## [7] "Acer sp." "Achillea millefolium"
## [9] "Achillea millefolium (yarrow)" "Achillea sp."
```

```
## Remove the special character

oba$Associated.plant <- str_replace_all(oba$Associated.plant, "\\x92", "")

## Remove the rows with no plant name.

oba <- oba[oba$Associated.plant != "",]

## Remove "net"

oba <- oba[oba$Associated.plant != "Net",]

## Fix yarrow
```

```

oba$Associated.plant[oba$Associated.plant == "Yarrow"] <- "Achillea millefolium"

## Remove a random weird one

oba <- oba[oba$Associated.plant != "Weedy yellow comp.",]

## Remove names in ()

oba$Associated.plant <- str_replace_all(oba$Associated.plant, "\\(..*?\\)", "")

## Still some issues with words after commas

oba$Associated.plant <- str_replace_all(oba$Associated.plant, ",.*$", "")

## Some have write space at the end of the string now

oba$Associated.plant <- str_replace_all(oba$Associated.plant, "\\s+$", "")

## And now there are a few instances where sp doesn't have a period.

oba$Associated.plant <- str_replace_all(oba$Associated.plant, " sp$", " sp.")

## Remove the or and everything after it, could also consider dropping these...

oba$Associated.plant <- str_replace_all(oba$Associated.plant, " or .*", "")

sort(unique(oba$Associated.plant))[1:10]

## [1] "Abelia sp." "Abronia latifolia" "Acer circinatum"
## [4] "Acer macrophyllum" "Acer palmatum" "Acer sp."
## [7] "Achillea millefolium" "Achillea sp." "Aclepias speciosa"
## [10] "Aesculus hippocastanum"

```

## Part 2: Making a column for plant resolution

- Some plant species are resolved to species/subspecies, others to genus and others to family. If there are two or three words, we can assume the plant is resolved to species and subspecies, respectively, except if the string ends in “sp.” If there is only one word, this could be a genus or a family name. Family names always end in “aceae”, for example Lamiaceae (mints), Asteraceae (daisies).

We want to make a new column called plantResolution and assign it to “Family”, “Genus” or “Species” depending on the level of resolution associated taxa is resolved to. We will do this in two steps.

First use regular expressions to count up the number of words in each element of associated taxa. Assign the count to a new column called plantTaxaWordCount. Print the first 50 elements.

Hint: `str_count` may be useful.

```

# plantResolution column, assign it to "Family", "Genus", "Species"

# use regular expressions to count up the number of words in each element of assoc taxa
# assign count to a new column called plantTaxaWordCount, print first 50

```

```
# str_count(string, pattern) Locate positions of pattern matches in a string

oba <- oba %>%
  mutate(plantTaxaWordCount = str_count(oba$Associated.plant, '\\b\\w+\\b'))

head(oba$plantTaxaWordCount, 50)

## [1] 1 1 1 1 1 2 2 2 2 2 2 2 1 1 1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1
## [39] 1 1 1 2 2 2 2 2 2 2 2
```

- b. Write a for loop to assigned each entry of the column plantResolution to be “family”, “genus” or “species”. `table()` the final result. Hint: Don’t forget to initialize the new column. Starting with all NAs may be useful. Hint hint: The function `ifelse` returns one value if a TRUE and another if FALSE. It could be useful depending on your approach. Hint hint hint: `grepl` will return TRUE or FALSE depending on whether it finds the pattern. Be careful with periods in patterns because alone they are a wild card character.

```
# writing for to assign each entry of plantResolution for family genus or species
# ifelse returns one value if a TRUE and another if FALSE
# grepl will return TRUE or FALSE depending on whether it finds the pattern

# initialize new column plantResolution
oba <- oba %>%
  add_column(plantResolution = NA)

for (i in 1:nrow(oba)) {

  if (oba$plantTaxaWordCount[i] == 2 | oba$plantTaxaWordCount[i] == 3){
    oba$plantResolution[i] <- "species"
  }

  else if (oba$plantTaxaWordCount[i] == 1 &&
    str_detect(oba$Associated.plant[i], "aceae\\b")){

    oba$plantResolution[i] <- "family"
  }

  else {oba$plantResolution[i] <- "genus"}

}

#unique(oba$plantResolution)
```

- c. For those that are identified to genus but are lacking an sp., add that now so that they will not be treated as separate plant species (i.e., Rosa vs Rosa sp.). You can do this with a regular expression and using `gsub` or `string_replace_all` or by counting up the number of words in `Associated.plant`.

```
oba$Associated.plant <- ifelse(
  (oba$plantResolution == "genus" & !str_detect(oba$Associated.plant, "sp.\\b")), str_replace_all(oba$A
  oba$Associated.plant
)
```

*## To check that it worked*

```
unique(oba$Associated.plant[oba$plantResolution == "genus"])
```

```
## [1] "Salix sp."
## [2] "Leucanthemum sp."
## [3] "Caprifoliaceae sp."
## [4] "Helenium sp."
## [5] "Composite sp."
## [6] "Magnoliopsida sp."
## [7] "Poales sp."
## [8] "Lathyrus japonicus var. maritimus sp."
## [9] "Cuscuta salina var. major sp."
## [10] "Crataegus sp."
## [11] "Cichorioideae sp."
## [12] "Boraginales sp."
## [13] "Convolvulus sp."
## [14] "Anthemideae sp."
## [15] "Cichorieae sp."
## [16] "Artemisia sp."
## [17] "Chamaemelum sp."
## [18] "Arctostaphylos sp."
## [19] "Phacelia sp."
## [20] "Magnoliophyta sp."
## [21] "Umbelliferae sp."
## [22] "Calochortus sp."
## [23] "Aguilegia sp."
## [24] "Asclepias sp."
## [25] "Sweep over Bellis perennis sp."
## [26] "Melilotus alba and thistle sp."
## [27] "Potentilla only plant in bloom sp."
## [28] "Poison hemlock and Steens Mtn thistle sp."
## [29] "Sweep in restoration area with Prunella vulgare sp."
## [30] "Tarweed sp."
## [31] "Cirsium sp."
```

- d. Create a new column called plantGenus that is the genus if the associated taxa was resolved to species or genus, and NA if it was resolved to family.

```
oba <- oba %>%
  add_column(plantGenus = NA)

for (i in 1:nrow(oba)) {

  if (oba$plantResolution[i] == "genus"){
    # takes the genus if it's == genus
    oba$plantGenus[i] <- oba$Associated.plant[i]
  }
}
```

```

else if(oba$plantResolution[i] == "species"){

  # takes the first part of species
  oba$plantGenus[i] <- str_extract(oba$Associated.plant[i], "[A-Za-z]+\b")

}

else {
  oba$plantGenus[i] <- NA
}
}

## To check finish with
table(oba$plantGenus)

```

```

##
##               Abelia
##                5
##             Abronia
##                2
##               Acer
##               40
##             Achillea
##              108
##             Aclepias
##               10
##             Aesculus
##               13
##             Agastache
##               50
##           Agulegia sp.
##                2
##               Alcea
##               18
##             Allium
##               67
##             Alyssum
##                1
##           Amelanchier
##               41
##             Amsinckia
##                2
##             Anaphalis
##                4
##             Anemone
##                1
##             Anethum
##                1
##             Angelica
##                4
##           Antennaria
##                9

```



##	Anthemideae sp.
##	11
##	Antirrhinum
##	6
##	Apocynum
##	32
##	Aquilegia
##	20
##	Arabis
##	2
##	Arbutus
##	1
##	Arctostaphylos
##	96
##	Arctostaphylos sp.
##	2
##	Arenaria
##	9
##	Arnica
##	51
##	Artemisia
##	15
##	Artemisia sp.
##	1
##	Asclepias
##	50
##	Asclepias sp.
##	19
##	Asparagus
##	7
##	Aster
##	75
##	Astragalus
##	22
##	Astragalus
##	1
##	Atriplex
##	1
##	Aurinia
##	1
##	Balsamorhiza
##	54
##	Barbarea
##	8
##	Bellardia
##	2
##	Bellis
##	214
##	Berberis
##	159
##	Beta
##	4
##	Betula
##	2

##	Bidens
##	4
##	Bistorta
##	15
##	Blepharipappus
##	9
##	Boechera
##	7
##	Boraginales sp.
##	6
##	Borago
##	2
##	Brassica
##	178
##	Brodiaea
##	12
##	Bromus
##	1
##	Calendula
##	9
##	Calochortus
##	43
##	Calochortus sp.
##	1
##	Calystegia
##	3
##	Camassia
##	102
##	Caprifoliaceae sp.
##	5
##	Cardaria
##	8
##	Carpenteria
##	1
##	Caryopteris
##	9
##	Castanea
##	2
##	Castilleja
##	1
##	Catalpa
##	1
##	Ceanothus
##	388
##	Centaurea
##	25
##	Centauria
##	1
##	Cerastium
##	19
##	Chaenactis
##	6
##	Chaenomeles
##	2

##	Chamaemelum sp.
##	1
##	Chamaenerion
##	51
##	Chorispora
##	20
##	Chrysolepis
##	3
##	Chrysothamnus
##	120
##	Chrysomanthus
##	5
##	Cichorieae sp.
##	5
##	Cichorioideae sp.
##	20
##	Cirsium
##	69
##	Cirsium sp.
##	17
##	Cistus
##	12
##	Clarkia
##	29
##	Claytonia
##	64
##	Cleome
##	13
##	Collinsia
##	12
##	Composite sp.
##	3
##	Convolvulus
##	4
##	Convolvulus sp.
##	3
##	Coreopsis
##	17
##	Coriandrum
##	1
##	Cornus
##	40
##	Cosmos
##	24
##	Cotinus
##	9
##	Crataegus
##	33
##	Crataegus sp.
##	9
##	Crepis
##	54
##	Crocidium
##	1

##	Crocus
##	2
##	Cryptantha
##	3
##	Cucurbita
##	3
##	Cuscuta salina var. major sp.
##	2
##	Cynoglossum
##	2
##	Cytisus
##	13
##	Dahlia
##	7
##	Damasonium
##	1
##	Dasiphora
##	23
##	Daucus
##	50
##	Delphinium
##	36
##	Descurainia
##	38
##	Deutzia
##	3
##	Dianthus
##	7
##	Dicentra
##	5
##	Dichelostemma
##	76
##	Digitalis
##	17
##	Dipsacus
##	21
##	Doronicum
##	20
##	Downingia
##	2
##	Drymocallis
##	6
##	Echinacea
##	23
##	Echinops
##	3
##	Elaeagnus
##	7
##	Epilobium
##	51
##	Ericameria
##	524
##	Erigeron
##	51

##	Eriodictyon
##	18
##	Eriogonum
##	1
##	Eriogonum
##	119
##	Eriophyllum
##	212
##	Erodium
##	12
##	Eruca
##	9
##	Erysimum
##	7
##	Erythronium
##	6
##	Escallonia
##	1
##	Escholtzia
##	3
##	Eschscholzia
##	279
##	Euonymus
##	8
##	Euphorbia
##	25
##	Fagopyrum
##	4
##	Foeniculum
##	19
##	Forsythia
##	12
##	Fragaria
##	76
##	Fraxinus
##	3
##	Gaillardia
##	26
##	Gentiana
##	1
##	Geranium
##	115
##	Geum
##	6
##	Gilia
##	39
##	Gnaphalium
##	1
##	Grindelia
##	79
##	Hackelia
##	37
##	Hastata
##	5

##	Hebe
##	5
##	Helenium
##	31
##	Helenium sp.
##	1
##	Helianthus
##	120
##	Heliopsis
##	9
##	Hemizonella
##	8
##	Heracleum
##	131
##	Hesperis
##	5
##	Heuchera
##	26
##	Hieracium
##	9
##	Hirschfeldia
##	1
##	Holodiscus
##	30
##	Horkelia
##	31
##	Humulus
##	1
##	Hyacinthoides
##	21
##	Hyacinthus
##	2
##	Hydrangea
##	12
##	Hydrophyllum
##	23
##	Hypericum
##	6
##	Hypochaeris
##	52
##	Hyssopus
##	4
##	Ilex
##	63
##	Iliamna
##	11
##	Iris
##	9
##	Isatis
##	2
##	Jacobaea
##	13
##	Jaumea
##	1

##	Kalmia
##	3
##	Lamium
##	9
##	Larkspur
##	3
##	Lasthenia
##	5
##	Lathyrus
##	10
##	Lathyrus japonicus var. maritimus sp.
##	2
##	Lavandula
##	20
##	Leontodon
##	5
##	Lepechinia
##	1
##	Lepidium
##	34
##	Leucanthemum
##	100
##	Leucanthemum sp.
##	15
##	Lewisia
##	2
##	Limnanthes
##	10
##	Linaria
##	2
##	Linum
##	6
##	Lithodora
##	1
##	Lithophragma
##	6
##	Lithospermum
##	28
##	Lobularia
##	1
##	Lomatium
##	64
##	Lonicera
##	66
##	Lotus
##	34
##	Lunaria
##	1
##	Lupinus
##	285
##	Lychnis
##	1
##	Madia
##	21

##	Magnoliophyta sp.	
##		7
##	Magnoliopsida sp.	
##		7
##	Mahonia	
##		16
##	Malus	
##		32
##	Malva	
##		1
##	Marah	
##		12
##	Matricaria	
##		4
##	Medicago	
##		10
##	Melelotus	
##		10
##	Melilotus	
##		37
##	Melilotus alba and thistle sp.	
##		4
##	Melissa	
##		2
##	Mentha	
##		13
##	Mentzelia	
##		4
##	Microseris	
##		19
##	Mimulus	
##		29
##	Monarda	
##		24
##	Monardella	
##		32
##	Myosotis	
##		21
##	Narcissus	
##		6
##	Nemophila	
##		19
##	Nepeta	
##		68
##	Oenanthë	
##		1
##	Oenothera	
##		4
##	Origanum	
##		30
##	Oxalis	
##		14
##	Packera	
##		1



##	Pastinaca
##	1
##	Penstemon
##	235
##	Perideridia
##	12
##	Perovskia
##	37
##	Petasites
##	1
##	Phacelia
##	299
##	Phacelia sp.
##	1
##	Philadelphus
##	20
##	Phlox
##	2
##	Photinia
##	3
##	Physocarpus
##	125
##	Pieris
##	6
##	Plagiobothrys
##	41
##	Plantago
##	1
##	Plectritis
##	100
##	Poales sp.
##	13
##	Poison hemlock and Steens Mtn thistle sp.
##	9
##	Polygonum
##	4
##	Potentilla
##	185
##	Potentilla only plant in bloom sp.
##	8
##	Prosartes
##	6
##	Prunella
##	9
##	Prunus
##	149
##	Pseudotsuga
##	15
##	Pseudoveronica
##	1
##	Purple
##	11
##	Purshia
##	31

##	Pyrus
##	24
##	Quercus
##	8
##	Ranunculus
##	54
##	Raphanus
##	194
##	Rheum
##	9
##	Rhododendron
##	88
##	Rhus
##	65
##	Ribes
##	300
##	Robinia
##	9
##	Romneya
##	1
##	Rosa
##	142
##	Rosmarinus
##	35
##	Rubus
##	473
##	Rudbeckia
##	14
##	Salix
##	66
##	Salix sp.
##	156
##	Salvia
##	24
##	Sanicula
##	6
##	Sarcobatus
##	3
##	Scabiosa
##	2
##	Scandix
##	61
##	Scilla
##	10
##	Scutellaria
##	4
##	Sedum
##	18
##	Senecio
##	75
##	Sidalacea
##	51
##	Sidalcea
##	113

##	Silene	
##		3
##	Silphium	
##		1
##	Sinapis	
##		17
##	Sisymbrium	
##		6
##	Sisyrinchium	
##		1
##	Small	
##		1
##	Solanum	
##		3
##	Solidago	
##		209
##	Sonchus	
##		19
##	Sorbus	
##		20
##	Sphaeralcea	
##		4
##	Sphenosciadium	
##		14
##	Spiraea	
##		172
##	Stachys	
##		7
##	Stephanomeria	
##		1
##	Styrax	
##		21
##	Sweep in restoration area with Prunella vulgare sp.	
##		1
##	Sweep over Bellis perennis sp.	
##		2
##	Symphoricarpos	
##		87
##	Symphyotrichum	
##		59
##	Symphytum	
##		8
##	Syringa	
##		17
##	Tanacetum	
##		25
##	Taraxacum	
##		244
##	Taraxia	
##		4
##	Tarweed sp.	
##		4
##	Tellima	
##		14

##	Teucrium
##	3
##	Thelypodium
##	135
##	Thermopsis
##	17
##	Thymus
##	1
##	Tithonia
##	2
##	Tonella
##	14
##	Toxicoscordion
##	6
##	Tradescantia
##	7
##	Tragopogon
##	15
##	Trifolium
##	111
##	Triteleia
##	9
##	Umbelliferae sp.
##	2
##	Vaccinium
##	58
##	Veratrum
##	2
##	Verbascum
##	3
##	Verbena
##	33
##	Veronica
##	42
##	Viburnum
##	13
##	Vicia
##	218
##	Viola
##	4
##	Weigela
##	6
##	Whipplea
##	51
##	Wisteria
##	21
##	Wyethia
##	88
##	Zinnia
##	1

Now you have nice clean plant data to make networks out of, or more easily count up the number of plant species in an area.