

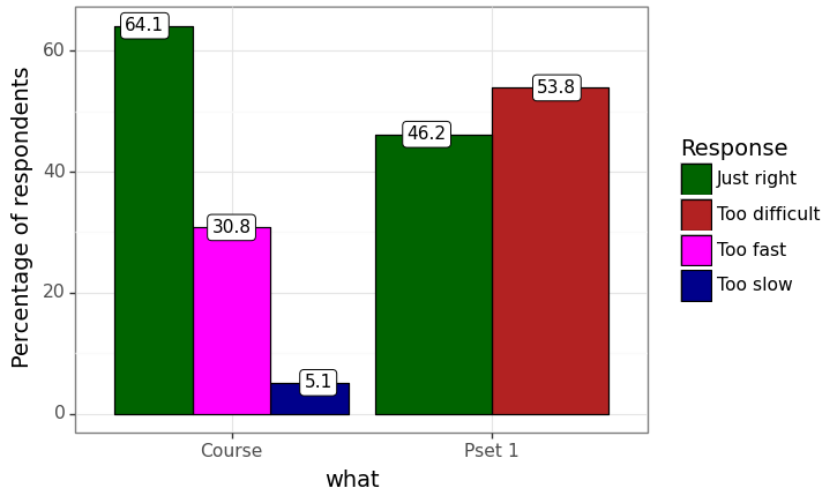
QSS20: Modern Statistical Computing

Unit 05: Merging (exact)

Goal for next few sessions

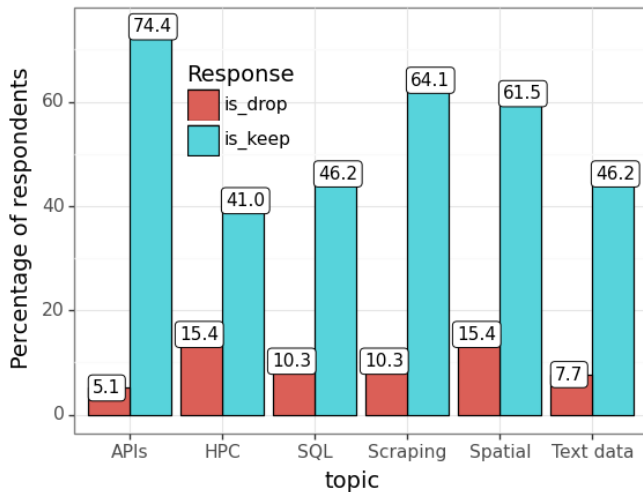
- ▶ **Today:** Exact merging: types of joins
 - ▶ Inner joins
 - ▶ Outer joins
 - ▶ Left joins
 - ▶ Right joins
- ▶ **Wednesday:** basic regex for two purposes:
 1. Clean join fields for exact matching/merges
 2. Clean join fields for fuzzy/probabilistic matching/merges
- ▶ **Next week:** fuzzy/probabilistic merging/record linkage

Summary of problem set 1 difficulty feedback

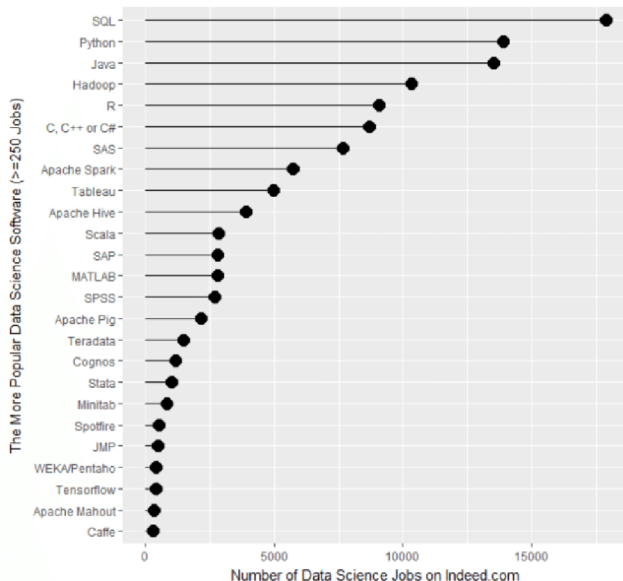


$$\frac{39 \text{ respondents}}{55 \text{ students}} = 71\% \text{ response rate}$$

Summary of problem set 1 topics feedback



Data science careers



Goals for today

- ▶ Practice git/GitHub from last week and review pset2 submission instructions
- ▶ Exact merging lecture
- ▶ Exact merging activity- see [03_merging_exact_blank.ipynb](#)

Goals for today

- ▶ Practice git/GitHub from last week and review pset2 submission instructions
- ▶ **Exact merging lecture**
- ▶ Exact merging activity- see [03_merging_exact_blank.ipynb](#)

Working example: have dataset on Dartmouth students and want to merge in background information about their district

► **Main or “left” dataset**

Student	Year	District	NCES ID
Rebecca	2021	New Trier High School	1728200
Jennifer	2022	Hanover High	3302670
Jason	2022	Homeschool	NA
⋮			

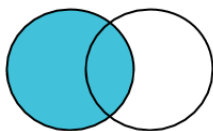
► **Auxiliary or “right” dataset**

District	NCES ID	% FRPL
New Trier HS	1728200	X%
Hanover HS	3302670	Y%
Lebanon HS	4107380	Z%
⋮		

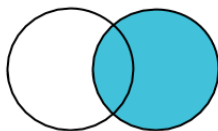
Possible join keys

- ▶ **Unique identifier:** used for “exact matching” — or a Yes/No match on that basis
 - ▶ E.g., is the NCES ID of New Trier found in the dataset of demographics?
- ▶ **Other identifiers:** can be used for either “exact match” or for “probabilistic/fuzzy matching”
 - ▶ **Probabilistic:** what’s the likelihood that “New Trier district” and “New Trier HS” are the same entity?

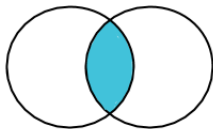
Conceptual overview of four types of joins



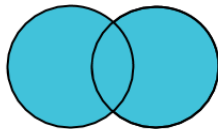
Left Join



Right Join



Inner Join



**Full Outer
Join**

Source: Trifacta

Inner join in this context

In words: “drop all students whose districts don’t appear in the demographics data; drop all districts that don’t appear in the Dartmouth student data”

► **Main or “left” dataset**

Student	Year	District	NCES ID
Rebecca	2021	New Trier High School	1728200
Jennifer	2022	Hanover High	3302670
Jason	2022	Homeschool	NA
⋮			

► **Auxiliary or “right” dataset**

District	NCES ID	% FRPL
New Trier HS	1728200	X%
Hanover HS	3302670	Y%
Lebanon HS	4107380	Z%
⋮		

Outer join in this context

In words: “keep all students from the student-level data; keep all schools from the school-level data; even if there’s not an overlap”

Student	Year	District	NCES ID	% FRPL
Rebecca	2021	New Trier High School	1728200	X%
Jennifer	2022	Hanover High	3302670	Y%
Jason	2022	Homeschool	NA	NA
NA	NA	NA	4107380	Z%
⋮				

Left join in this context

In words: “keep all students from the student-level data; drop any school from the school-level data that doesn’t merge onto a student”

► **Main or “left” dataset**

Student	Year	District	NCES ID
Rebecca	2021	New Trier High School	1728200
Jennifer	2022	Hanover High	3302670
Jason	2022	Homeschool	NA
⋮			

► **Auxiliary or “right” dataset**

District	NCES ID	% FRPL
New Trier HS	1728200	X%
Hanover HS	3302670	Y%
Lebanon HS	4107380	Z%
⋮		

Right join in this context

In words: “drop students who don’t have a school in the school-level data; keep all schools from the student-level data even those that don’t merge onto any student”

► **Main or “left” dataset**

Student	Year	District	NCES ID
Rebecca	2021	New Trier High School	1728200
Jennifer	2022	Hanover High	3302670
Jason	2022	Homeschool	NA
⋮			

► **Auxiliary or “right” dataset**

District	NCES ID	% FRPL
New Trier HS	1728200	X%
Hanover HS	3302670	Y%
Lebanon HS	4107380	Z%
⋮		

DataCamp versus slide syntax

- ▶ DataCamp modules generally used this syntax for merges:

```
1 merged_df = df1.merge(df2 ,  
2                       how = "left" ,  
3                       on = "something" )
```

- ▶ Slides/solution code will tend to use this syntax:

```
1 merged_df = pd.merge(df1 , df2 ,  
2                       how = "left" ,  
3                       on = "something" )
```

- ▶ They produce identical answers so use whichever comes more naturally (I use latter because it's more similar to base R syntax)
- ▶ In addition, feel free to use self joins if useful but we won't be focusing a lot on those

How do we code these different types of joins in practice?
Example with left join and join key has same colname in both

```
1  
2 ## perform a left join on the student data  
3 ## and schools data  
4 stud_wschoool = pd.merge(students ,  
5                           schools ,  
6                           how = "left" ,  
7                           on = "NCES ID" ,  
8                           indicator = "student_mergestatus" )
```

- ▶ **how**: argument to tell it inner, left, right, outer, or cross; defaults to inner
- ▶ **on**: name of join key (in this case single key)
- ▶ **indicator**: optional arg to add a col to the resulting data (string is what to call it) that helps diagnose merge status; good for post-merge dx

Example with inner join and join key has different name

```
1  
2 ## perform a left join on the student data  
3 ## and schools data  
4 stud_wschoool = pd.merge(students ,  
5                         schools ,  
6                         how = "inner" ,  
7                         left_on = "NCES ID" ,  
8                         right_on = "ncesnumeric")
```

Example with left join and multiple join keys

```
1  
2 ### perform a left join on the student data  
3 ### and schools data  
4 stud_wschoool = pd.merge(students ,  
5                           schools ,  
6                           how = "left" ,  
7                           left_on = ["NCES ID" ,  
8                                       "Dist name" ] ,  
9                           right_on = ["ncesnumeric" ,  
10                                      "distnamechar" ] ,  
11                          indicator = "student_mergestatus" )
```

Example with left join, multiple join keys, and some overlapping, non-join columns that we want to differentiate

```
1  
2 ### perform a left join on the student data  
3 ### and schools data  
4 stud_wschoool = pd.merge(students ,  
5                           schools ,  
6                           how = "left" ,  
7                           left_on = ["NCES ID" ,  
8                                       "Dist name" ] ,  
9                           right_on = ["ncesnumeric" ,  
10                                      "distnamechar" ] ,  
11                           indicator = "student_mergestatus" ,  
12                           suffixes = ["_students" ,  
13                                       "_schools" ] )
```

Non-exhaustive checklist of merge diagnostics

1. How many rows were in each data before the merge? What about after?
2. If doing a left join, did we properly retain all left-hand side rows?
3. **For strings as join keys:** if a lot of rows were lost in a merge, could that be due to spelling/punctuation variations in a character join key?
4. **For numeric identifiers as join keys:** if a lot of rows were lost in a merge, could that be due to things like the id having leading or lagging zeros and those being stripped at some stage? (e.g., one dataset identifies an entity as 002548; another as 2548)

Next up: basic regex to improve match rates for strings as join keys

- In example below, what if we didn't have the NCES ID numeric identifier? Ways to improve match rates for spelling variations (sometimes called `entity resolution`)

Student	Year	District
Rebecca	2021	New Trier High School
Jennifer	2022	Hanover High
Jason	2022	Homeschool
⋮		

District	% FRPL
New Trier HS	X%
Hanover HS	Y%
Lebanon HS	Z%
⋮	

Goals for today

- ▶ Practice git/GitHub from last week and review pset2 submission instructions
- ▶ Exact merging lecture
- ▶ **Exact merging activity-** see [03_merging_exact_blank.ipynb](#)

- ▶ `public_data/sd_df.csv`: sample of business tax certificates for San Diego-based businesses— each row represents one unique business; cols for industry (6-digit NAICS code)

account_key	dba_name	council_district	naics_code	naics_description	naics_
1974000448	ERNST & YOUNG LLP	cd_1	541211	OFFICES OF CERTIFIED PUBLIC ACCOUNTANTS	
1974011093	HECHT SOLBERG ROBINSON GOLDBERG & BAGLEY LLP	cd_3	5411	LEGAL SERVICES	
1978039819	RSM US LLP	cd_1	541211	OFFICES OF CERTIFIED PUBLIC ACCOUNTANTS	
1978042092	THORSNES BARTOLOTTA MCGUIRE LLP	cd_3	5411	LEGAL SERVICES	
1979046817	KORENIC & WOJDOWSKI LLP	cd_7	5412	ACCOUNTING/TAX PREP/BOOKKEEP/PAYROLL SERVICES	

- ▶ `public_data/naics_df.csv`: exhaustive listing of all 6-digit NAICS codes from the Census Bureau with added information

naics	naics_description
111140	Wheat Farming
111160	Rice Farming
111150	Corn Farming
111110	Soybean Farming
111120	Oilseed (except Soybean) Farming

- ▶ **General goal:** match the two to investigate things like which industries are *not* represented in the San Diego small businesses

Outline of notebooks

- ▶ Blank notebook to work on with partners/small group: [03_merging_exact_blank.ipynb](#)
- ▶ Input data: `public_data/sd_df.csv` and `public_data/naics_df.csv`
- ▶ Solutions notebook to review after class (better to avoid looking during activity!): [03_merging_exact_solutions.ipynb](#)
- ▶ In case of interest, notebook where I cleaned/prepped the data to make easier to analyze: [03helper_merging_dataprep.ipynb](#)