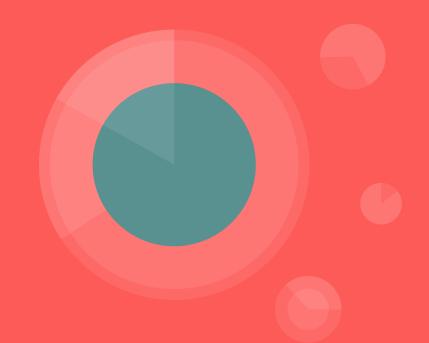
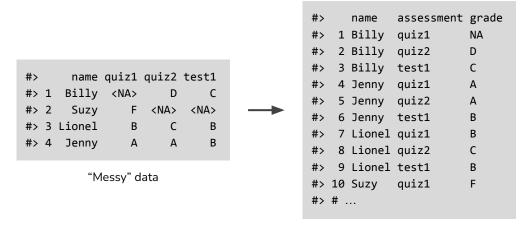
# Database normalization



### What is normalization?

- Re-organizing datasets to emphasize relationships and modularity
  - There are increasingly strict levels of normalization ("normal forms")
- "Tidy" data is a form of normalization (similar to 3rd normal form)



Normalized/"tidied" data

### Why normalize?

Protects against errors and makes data entry easier
Employees' Skills

Employee ID	Employee Address	Skill				
426	87 Sycamore Grove	Typing				
426	87 Sycamore Grove	Shorthand				
519	94 Chestnut Street	Public Speaking				
519	96 Walnut Avenue	Carpentry				

← A less-normalized setup can make contradictions easier (e.g. conflicting address info).

- Reduces the need to refactor databases as more data is added. (What if we wanted to add ZIP codes?)
- Simplifies queries
  - Less shuffling columns, manual parsing, or wrangling nested data
  - Consistency makes statistical analyses more straightforward



#### What does this have to do with INFO 523?

- Essential data wrangling step for more complex datasets (tidyverse needs "tidy" data == normalization)
- Often a goal of reshaping or organizing data relationally (see Kabacoff (2015), Ch. 5; Wickham and Grolemund (2016), Ch. 12-13)
- Normalization's scalability makes it a very good idea for "big" data (data stores, data warehousing, etc.)



### Real-world applications

- Normalization is a fundamental principle of "relational database management systems" (RMDBS): MySQL, Postgres, Oracle Database, Apache Hive, etc.
- Proper use of databases with these systems requires an understanding of database normalization.

## Let's try it!

### Fundamental ideas

- Sometimes data should might work better in multiple tables!
  - The more narrowly a table describes a relationship, the more "normalized"
- Rows have "primary keys" that identify them, and then additional columns
  - You get to pick which columns are "primary keys"
  - Primary keys are used to cross-reference different tables
  - Columns ("non-primary keys") may or may not "depend" on the primary keys—this is what they mean by "relational" data
  - The exact terminology is very particular, I'll stick to "(primary) keys" and "columns"

### Normalization levels (for reference)

	<b>UNF</b> (1970)	<b>1NF</b> (1970)	<b>2NF</b> (1971)	<b>3NF</b> (1971)	<b>EKNF</b> (1982)	<b>BCNF</b> (1974)	<b>4NF</b> (1977)	ETNF (2012)	<b>5NF</b> (1979)	<b>DKNF</b> (1981)	<b>6NF</b> (2003)
Primary key (no duplicate tuples) <sup>[5]</sup>	V	V	V	V	V	V	V	V	V	V	V
Atomic columns (cells cannot have tables as values) <sup>[6]</sup>	×	V	V	V	V	V	V	V	V	V	V
No partial functional dependencies of non-prime attributes on candidate keys) <sup>[6]</sup>	X	×	V	V	V	V	V	V	V	V	V
No transitive functional dependencies of non-prime attributes on candidate keys) <sup>[6]</sup>	×	×	×	V	V	<b>V</b>	V	V	V	V	V
Every non-trivial functional dependency either begins with a superkey or ends with an elementary prime attribute	×	×	×	×	V	V	V	V	V	V	N/A
Every non-trivial functional dependency begins with a superkey	×	×	×	×	×	V	V	V	V	V	N/A
Every non-trivial multivalued dependency begins with a superkey	×	×	×	×	×	×	V	V	V	V	N/A
Every join dependency has a superkey component [9]	×	×	×	×	×	×	×	V	V	V	N/A
Every join dependency has only superkey components	×	×	×	×	×	×	×	×	V	V	N/A
Every constraint is a consequence of domain constraints and key constraints	×	×	×	×	X	X	X	X	X	V	×
Every join dependency is trivial	X	X	X	X	X	X	X	X	X	X	V

### Conclusion

### **Conclusions**

- You may never need to use 6NF, but normalization is handy to think about
  - (And you've probably already used it!)
- Be careful of your resources—if you want to know more, read up on the original papers
  - Older: Edgar F. Codd
  - Newer: Christopher J Date