



190F Foundations of Data Science

Spring 2020

Lecture 4

Building Tables
& Census

Announcements

Tables Review

Table Structure

- A Table is a sequence of labeled columns
- Labels are strings
- Columns are arrays, all with the same length

The diagram illustrates a table structure with three columns: Name, Code, and Area (m2). The first two columns are highlighted with a blue rounded rectangle, and the last two columns are highlighted with a blue rounded rectangle. A blue callout box labeled 'Label' points to the 'Code' header. A blue callout box labeled 'Row' points to the 'Nevada' row. A blue callout box labeled 'Column' points to the 'Code' column.

Name	Code	Area (m2)
California	CA	163696
Nevada	NV	110567

Ways to create a table

- `Table.read_table(filename)` - reads a table from a spreadsheet
- `Table()` - an empty table
- and...

Arrays → Tables

- `Table().with_column(label, data)` - creates a table with a single column; `data` is an array
 - `Table().with_columns(label1, data1, ...)` - creates a table, with an array of data for each column
-

Table Methods

- Creating and extending tables:
 - `Table().with_columns` and `Table.read_table`
 - Finding the size: `num_rows` and `num_columns`
 - Referring to columns: labels, relabeling, and indices
 - `labels` and `reabeled`; column indices start at 0
 - Accessing data in a column
 - `column` takes a label or index and returns an array
 - Using array methods to work with data in columns
 - `item`, `sum`, `min`, `max`, and so on
 - Creating new tables containing some of the original columns:
 - `select`, `drop`
-

Examples

The table `students` has columns `Name`, `ID`, and `Score`.
Write one line of code that evaluates to:

a) A table consisting of only the column labeled `Name`

```
students.select('Name')
```

b) The largest score

```
students.column('Score').max()  
max(students.column('Score'))
```

Minard's Map

Charles Joseph Minard, 1781-1870

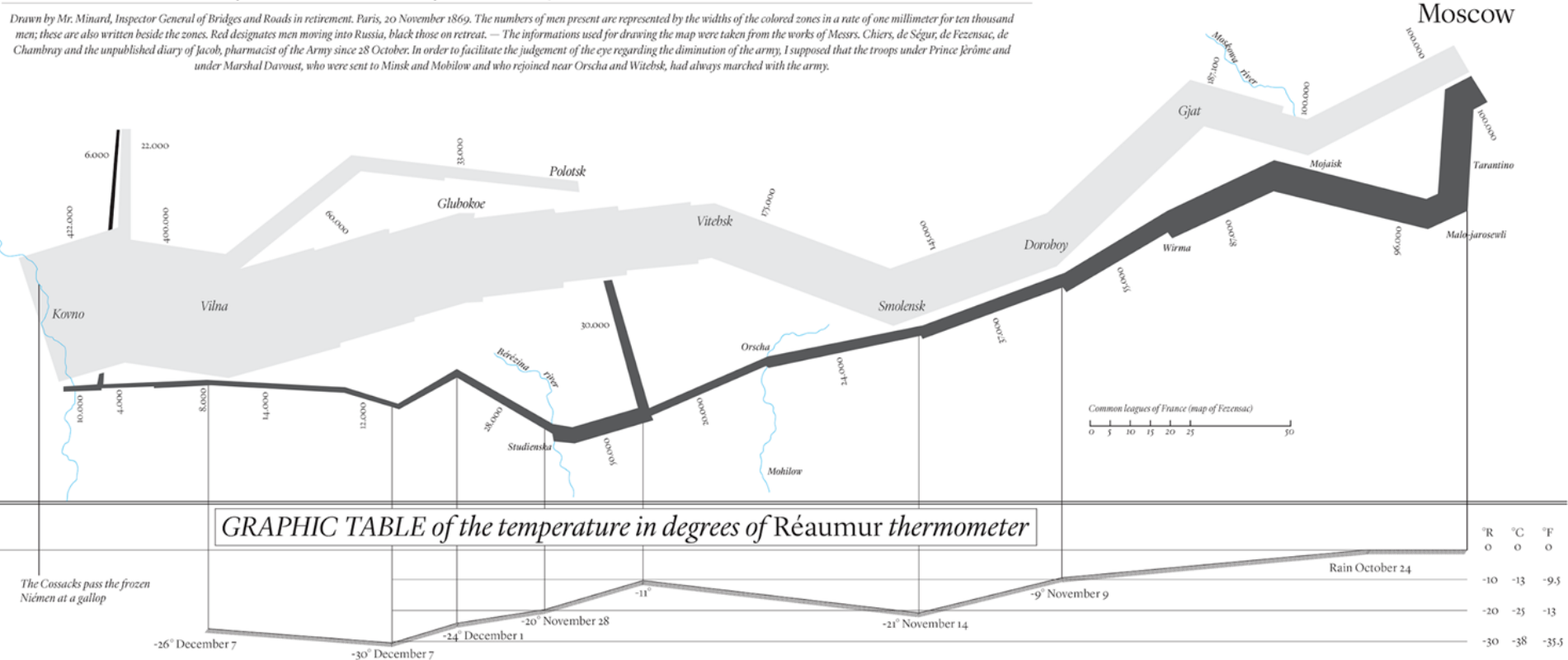


- French civil engineer who created one of the greatest graphs of all time
- Visualized Napoleon's 1812 invasion of Russia, including
 - the number of soldiers
 - the direction of the march
 - the latitude and longitude of each city
 - the temperature on the return journey
 - Dates in November and December

Visualization of 1812 March

FIGURATIVE MAP of the successive losses in men of the French Army in the RUSSIAN CAMPAIGN OF 1812-1813

Drawn by Mr. Minard, Inspector General of Bridges and Roads in retirement, Paris, 20 November 1869. The numbers of men present are represented by the widths of the colored zones in a rate of one millimeter for ten thousand men; these are also written beside the zones. Red designates men moving into Russia, black those on retreat. — The informations used for drawing the map were taken from the works of Messrs. Chiers, de Ségur, de Fezensac, de Chambray and the unpublished diary of Jacob, pharmacist of the Army since 28 October. In order to facilitate the judgement of the eye regarding the diminution of the army, I supposed that the troops under Prince Jérôme and under Marshal Davoust, who were sent to Minsk and Mohilow and who rejoined near Orscha and Witebsk, had always marched with the army.



Different types of data

float:
decimal number

Longitude	Latitude	City	Direction	Survivors
32	54.8	Smolensk	Advance	145000
33.2	54.9	Dorogobouge	Advance	140000
34.4	55.5	Chjat	Advance	127100
37.6	55.8	Moscou	Advance	100000
34.3	55.2	Wixma	Retreat	55000
32	54.6	Smolensk	Retreat	24000
30.4	54.4	Orscha	Retreat	20000
26.8	54.3	Moiodexno	Retreat	12000

string:
text

int:
integer

Lists

Lists are Generic Sequences

A list is a sequence of values (just like an array), but the values can all have different types

```
[2+3, 'four', Table().with_column('K', [3, 4])]
```

If you create a table column from a list, it will be converted to an array automatically

(Demo)

Take

Take Rows, Select Columns

The `select` method returns a table with only some columns

The `take` method returns a table with only some rows

- Rows are numbered, starting at 0
- Taking a single number returns a one-row table
- Taking a list of numbers returns a table as well

(Demo)

The where method

- `t.where(label, condition)` - constructs a new table with just the rows that match the condition

(Demo)

Manipulating Rows

- `t.sort(column)` sorts the rows in increasing order
 - `t.take(row_numbers)` keeps the numbered rows
 - Each `row` has an index, starting at 0
 - `t.where(column, are.condition)` keeps all rows for which a column's value satisfies a condition
 - `t.where(column, value)` keeps all rows containing a certain value in a column
 - `t.with_row` makes a new table that has another row
-

Discussion Questions

The table `nba` has columns `NAME`, `POSITION`, and `SALARY`.

- a) Create an array containing the names of all point guards (PG) who make more than \$15M/year

```
nba.where(1, 'PG').where(2, are.above(15)).column(0)
```

- b) After evaluating these two expressions in order, what's the result of the second one?

```
nba.with_row(['Samosa', 'Mascot', 100])  
nba.where('NAME', are.containing('Samo'))
```

Census Data

The Decennial Census

- Every ten years, the Census Bureau counts how many people there are in the U.S.
 - In between censuses, the Bureau estimates how many people there are each year.
 - Article 1, Section 2 of the Constitution:
 - “Representatives and direct Taxes shall be apportioned among the several States ... according to their respective Numbers ...”
-

Analyzing Census Data

Leads to the discovery of interesting features and trends in the population

(Demo)

Census Table Description

- Values have column-dependent interpretations
 - The SEX column: 1 is *Male*, 2 is *Female*
 - The POPESTIMATE2010 column: *7/1/2010 estimate*
- In this table, some rows are sums of other rows
 - The SEX column: 0 is *Total* (of *Male* + *Female*)
 - The AGE column: 999 is *Total* of all ages
- Numeric codes are often used for storage efficiency
- Values in a column have the same type, but are not necessarily comparable (AGE 12 vs AGE 999)