

# Data100 Sp22 Disc 4

## Regex/Viz!



**Attendance:**

**<https://tinyurl.com/disc4michelle>**

# Announcements

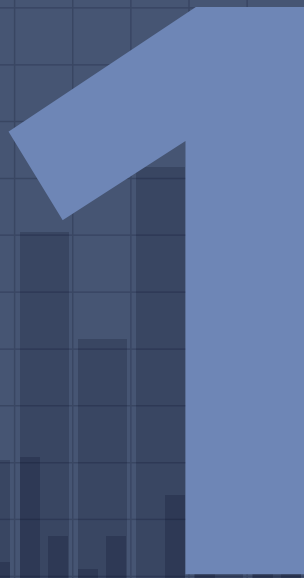
## Due Dates

- Homework 4 due Feb 17
- Lab 4 due Feb 15
- Weekly Check 4 due Feb 14

## Other

- Review session Feb19-20! (Weekend, will be recorded)
- Midterm Feb. 24

# Regex



# Basic Regex

operation	example	matches	does not match
or	TEA BOBA	TEA BOBA	every other string
character class	[A-Za-z]*	helloWorld BoBaTEa	b0b4T AB_AB
beginning of line (carat)	^word.*	word wordle	sword
end of line	.*word\$	word sword	wordle
character class negation (carat within [])	[^a-z]+	PEPPERS3982 17211!↑å	porch CLAmS

# Repeating Characters

operation	example	matches	does not match
Zero or more	BOB*A	BOA BOBBBBBA	BOOA BBBA
One or more	BOB+A	BOBA BOBBBBBA	BOA BBBA
zero or one	cats?	cat cats	any other string
repeated exactly {a} times	c[ao]{2}t	coat caat	caocat cat
repeated from a to b times: {a,b}	c[ao]{1,2}t	coat cat	caocat ct
Parenthesis (also used in catching groups)	(CA)+T	CACACAT CAT	CAAT T

# Character Groups

1. **.** (**period**): any character apart from a '\n' (newline)
2. **\w**: a "word" character (letter, digit, or underscore) [a-zA-Z0-9\_]
3. **\s**: matches a single whitespace character
4. **\d**: decimal digit [0-9]
5. **\**: escapes a special regex character. e.g. \( means match a opening parenthesis (

# Regex in Python

1. `re.match(pattern, string, ...)`: match if zero or more characters at the beginning of string match given pattern
2. `re.search(pattern, string, ...)`: scan through string looking for the first location where the regular expression pattern produces a match
3. `re.findall(pattern, string, ...)`: returns all non-overlapping matches of pattern in string as a list of strings
4. `re.sub(pattern, repl, string, ...)`: return the string obtained by replacing the non-overlapping occurrences of pattern in string by the replacement `repl`

# Other features

## 1. Greedy vs Lazy

- a. Greedy: match as many times as possible
- b. Lazy: match as few times as possible
- c. `+?` (lazy): `a+?` only ever matches a max of 1 `a`

## 2. Order of operations

ERE Precedence (from high to low)		
1	Collation-related bracket symbols	<code>[==]</code> <code>[::]</code> <code>[..]</code>
2	Escaped characters	<code>\&lt;special character&gt;</code>
3	Bracket expression	<code>[]</code>
4	Grouping	<code>()</code>
5	Single-character-ERE duplication	<code>*</code> <code>+</code> <code>?</code> <code>{m,n}</code>
6	Concatenation	
7	Anchoring	<code>^</code> <code>\$</code>
8	Alternation	<code> </code>



# Quickfire Regex

Which of these regex strings:

Matches: jossh, jossssh

Doesn't match: josh, joh, bjosssh

1. `jos[s]+h`
2. `^jos[s]+h`
3. `^jos[s]*h`
4. `jos[s]*h`

# Quickfire Regex

Which of these regex strings:

Matches: jossh, jossssh

Doesn't match: josh, joh, bjosssh

1. `jos[s]+h`: matches b(josssh)
2. `^jos[s]+h`
3. `^jos[s]*h`: matches josh
4. `jos[s]*h`: matches b(josssh)

# Quickfire Regex

Which strings would this Regex match?

```
[\\w]?[\\d]{2,3}[\\w]+[A-Z]
```

1. j05hl15A
2. j0shli5A
3. 05lisA
4. j055shhhli5a

# Quickfire Regex

Which strings would this Regex match?

`[\w]?[\d]{2,3}[\w]+[A-Z]`

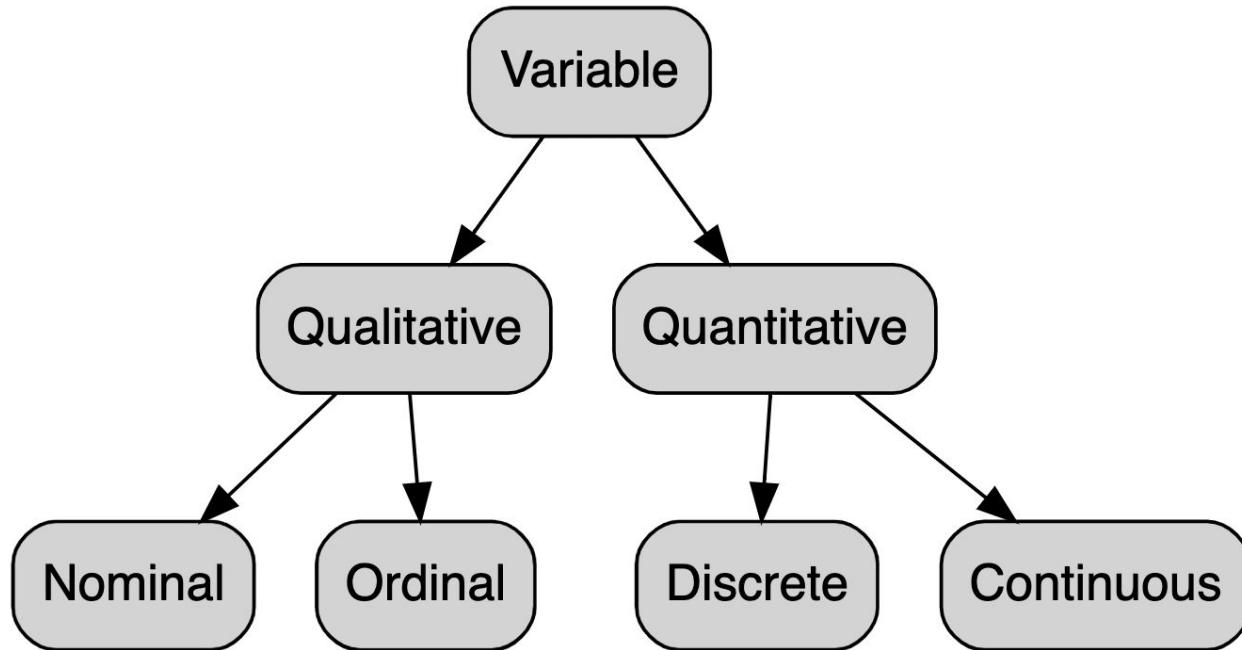
1. `j05hl15A`
2. `j0shli5A` : at least 2 digits needed for `[\d]{2,3}`
3. `05lisA`
4. `J055shhhli5A`: needs to end with an uppercase letter

# Visualizations

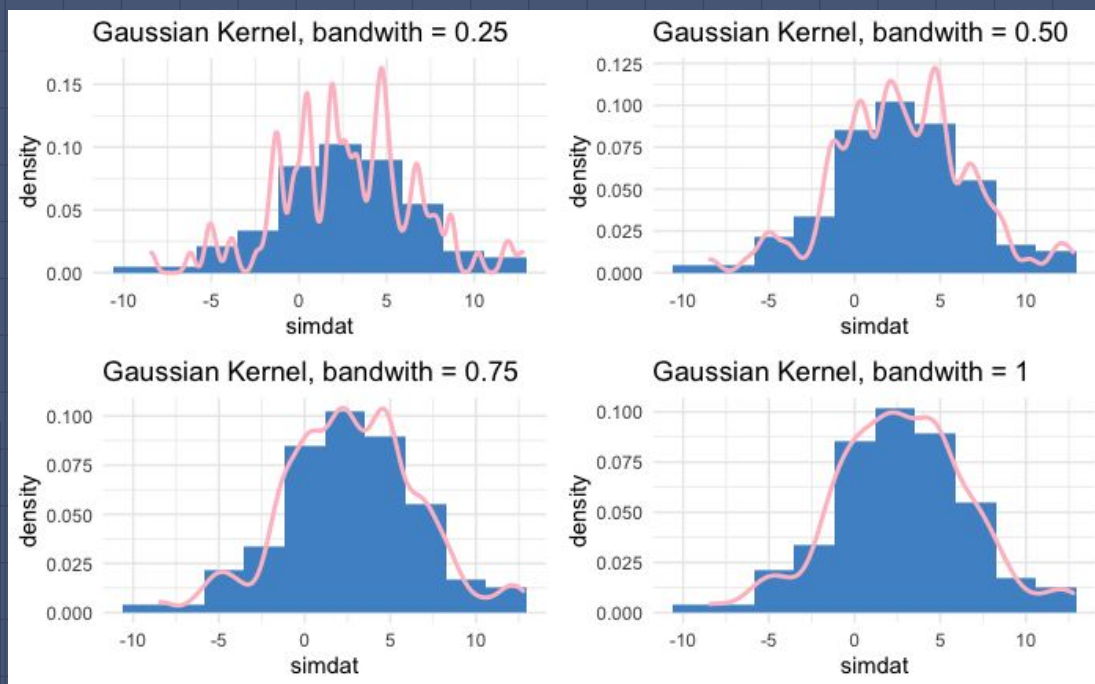
A large, light blue number 2 is positioned on the right side of the image. The background is a dark blue grid with a faint silhouette of a bar chart at the bottom. The number 2 is a simple, bold, sans-serif font.

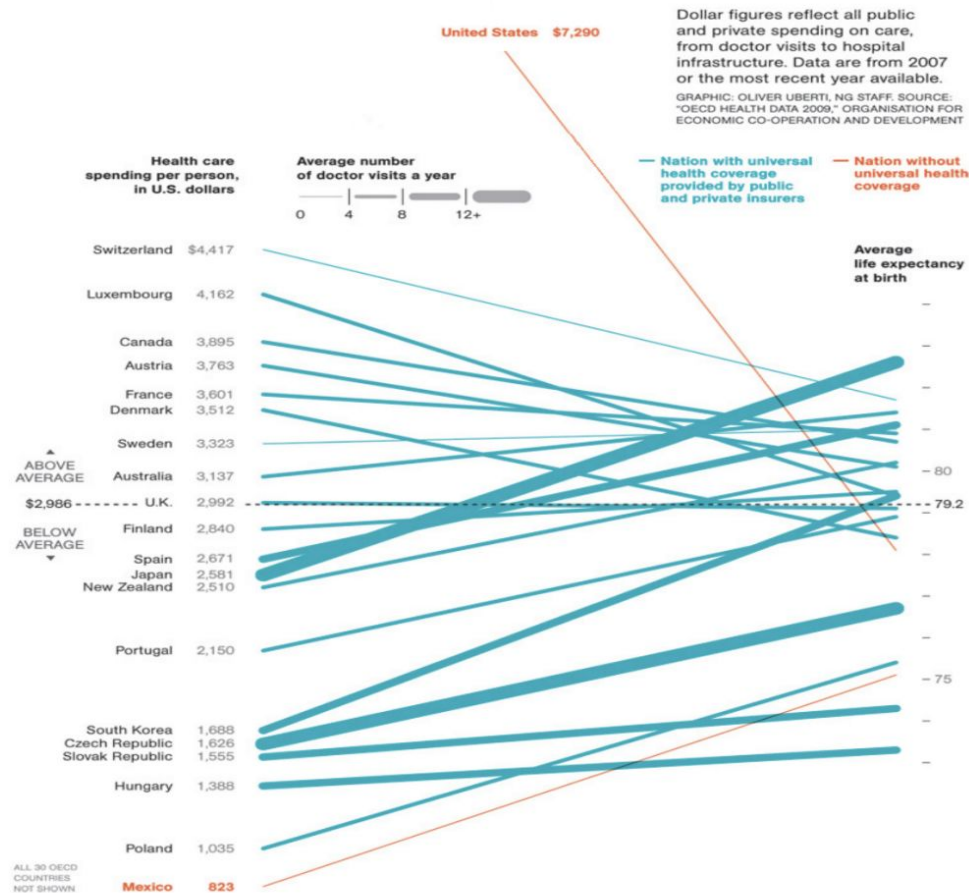
2

# Types of Variables



# KDEs







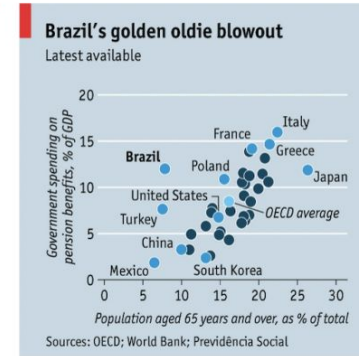
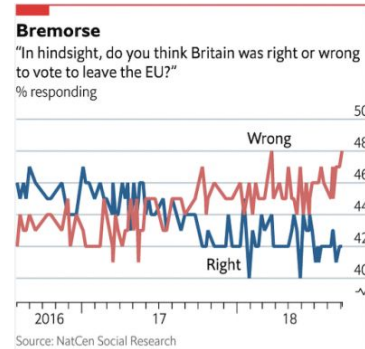
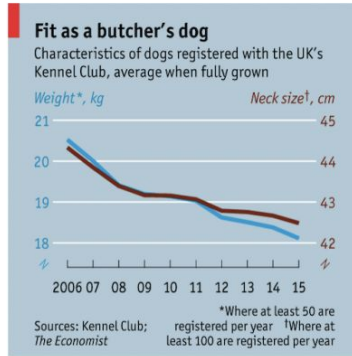
(a) Five variables are being represented visually in this graphic. What are they and what are their feature types (ie qualitative, quantitative, nominal, ordinal)?

(c) How can we figure out how to interpret the visual qualities of the plot, e.g., how do we know what a color represents?

(e) Make 3 observations about the figure. Describe the feature that you are basing your observation on.

For example, South Korea's expenditure on health care is comparable to Eastern European countries (and among the lowest of all countries plotted), but the life expectancy is much higher than the Eastern European countries. In the plot we see that the left endpoint of South Korea's line segment is near the Eastern European countries, but the slope of the line segment is much steeper.

9. Creating visualizations that represent data accurately and that support the narrative we wish to create is no easy task. Even the journalists and editors at *The Economist*, a newspaper known for its compelling, data-driven articles, have been known to make blunders. Three of their ill-thought-out plots are presented below. Consider what aspects of the visualizations are misleading, and think of ways in which you can remedy them.



*Hint:* The datapoints in the rightmost plot are shaded based on whether or not they are labeled.